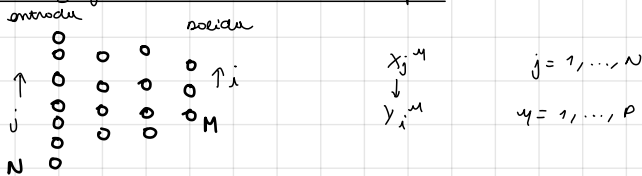


Aprendizaje en redes multicapa



Sistema lineal

$$g(x) = x$$

función de activación de una neurona

$$O_i^m = \sum_{j=1}^N w_{ij} x_j^m$$

salida real

Puedemos encontrar tales que $O_i^m = y_i^m$

$$w_{ij} ()$$

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

$$w_{ij} = \frac{1}{N} \sum_{m=1}^P y_i^m (Q^{-1})_{mv} x_j^m$$

$$Q_{mv} = \frac{1}{N} \sum_j x_j^m x_j^v \quad (P \times P)$$

es el que cumple con $O_i^m = y_i^m$

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

Aprendizaje por gradiente

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

$$E = 0 \Leftrightarrow y_i^m = O_i^m \quad i=1, \dots, M \quad m=1, \dots, P$$

Resolver el problema de aprendizaje \Leftrightarrow mínimo absoluto de E

$$\bar{w}_0$$

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

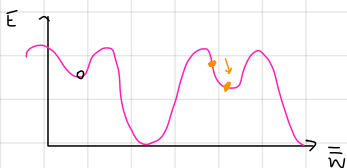
$$\Rightarrow E(\bar{w}_0 + \Delta \bar{w}) = E(\bar{w}) + \overline{\nabla} E(\bar{w}) \cdot \Delta \bar{w} + \mathcal{O}(\Delta \bar{w}^2)$$

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

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

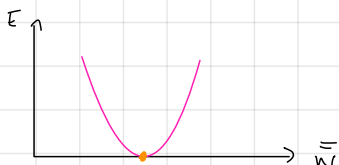
para asegurar esto pedimos η pequeño

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

Colocemos de manera explícita $E(\bar{w})$ para sist lineales ($g(x) = x$)

$$E(\bar{w}) = \frac{1}{2} \sum_{i=1}^n (y_i - \sum_j w_{ij} x_j)^2 = \dots + (y_k - \sum_j w_{kj} x_j)^2 + \dots$$

$$\frac{\partial E}{\partial w_{ke}} = - \sum_i (y_i - \sum_j w_{ij} x_j) x_{ie}$$

solo el término con $i=k$ sobrevive a w_{ke}

δe^k término de error

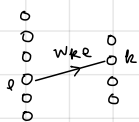
no continuos (mo 0,1 como en el perceptron)

$$\Delta w_{ke} = \eta \sum_{i=1}^p \delta_k^i x_{ie}$$

$\delta_{ke} = y_k - o_k$

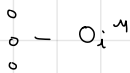
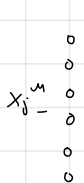
introduce regla de Adeline

diff. sólido deseado y \rightarrow error en la salida



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

una capa, no lineal (entrada - salida)



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

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

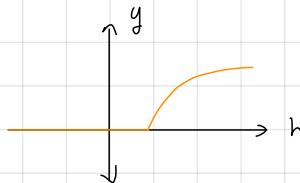
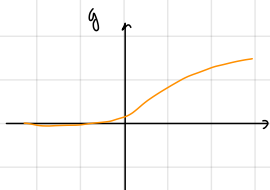
$$\frac{\partial E}{\partial w_{ke}} = \sum_u \left(y_k^u - g \left(\underbrace{\sum_{j=1}^N w_{kj} x_j^u}_{\substack{O_k^u \\ \delta_k^u}} \right) \right) g' \left(\underbrace{\sum_{j=1}^N w_{kj} x_j^u}_{h_k^u} \right) x_e^u$$

h_k^u
 entrada a neuronas en la última capa

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

posibles g :

(H4
integrable
en $\pm \infty$)



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

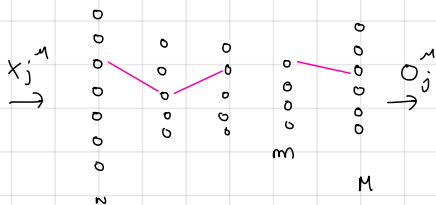
Podríamos definir la esperanza de otra manera:

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

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

se redefinió δ pero la forma funcional es la misma

Curso general: redes multicapa no-lineales



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

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

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

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

regla de la cadena

co/nea los
gradientes
para todos los conexiones a la vez
impracticable!

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

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

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

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

salida de la capa

$$E(\{O_i^m\}) = E\left(g\left(\sum_{h,i} w_{ih} O_h^m\right)\right), \frac{\partial E}{\partial w_{kl}} = \frac{\partial E}{\partial h_k^m} \frac{\partial h_k^m}{\partial w_{kl}} = \frac{\partial E}{\partial h_k^m} g'(h_k^m) O_l^m$$

- δ_k^m
mismo
que
el
indicio
de esta
neurona

al comenzar la salida, se comienza el entreno de las neuronas de la capa siguiente m .

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

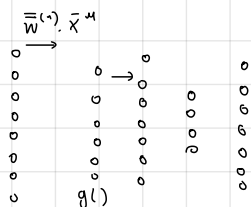
$\rightarrow \frac{\partial h_m^m}{\partial h_k^m} = w_{mk} g'(h_k^m)$

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

↓ 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

- aplico los w a los entornos
- calculo $w_1 x$
- aplico $g_1 \rightarrow g(w_1 x)$
- calculo $w_2 g_1(w_1 x)$
- aplico $g_2 \rightarrow g_2(w_2 g_1(w_1 x))$
- ⋮

Para la net capa

$$\frac{\partial E}{\partial h_m^m} = - \sum_k (y_k^m - g(h_k^m)) g'(h_k^m)$$

$$\frac{\partial E}{\partial h_{m-1}^m} = \frac{\partial E}{\partial h_m^m} \frac{\partial h_m^m}{\partial h_{m-1}^m}$$

$$= \frac{\partial E}{\partial h_m^m} w_{m-1, m} g'(h_{m-1}^m)$$

$$\frac{\partial E}{\partial h_{m-2}^m} = \frac{\partial E}{\partial h_{m-1}^m} \frac{\partial h_{m-1}^m}{\partial h_{m-2}^m} = \frac{\partial E}{\partial h_{m-1}^m} w_{m-2, m-1} g'(h_{m-2}^m)$$

- luego a la última capa y obtenemos los salidos \bar{o}
- calculo $\bar{y} - \bar{o}$
- luego prod pto a pto con la $g'(\bar{h}) : (\bar{y} - \bar{o}) \otimes g'(\bar{h})$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \cdot \begin{pmatrix} e & f \\ g & h \end{pmatrix} = \begin{pmatrix} ae & bf \\ cg & dh \end{pmatrix}$$

Queremos minimizar la función de costo E



$$\frac{\partial E}{\partial w_{jk}^e} = \frac{\partial E}{\partial h_j^e} \frac{\partial h_j^e}{\partial w_{jk}^e}$$

neutroa en capa $e-1$

neutroa en capa e

con

$$h_j^e = \sum_{k=1}^m w_{jk}^e o_k^{e-1} + b_j^e$$

m : # de neuronas en la capa $e-1$

($\frac{\partial E}{\partial w_{jk}^e}$: como varia E al cambiar una conexión $j-k$ en capa e

$\frac{\partial E}{\partial h_j^e}$ como varia E con el input de la neurona j en capa e

$\frac{\partial h_j^e}{\partial w_{jk}^e}$ como varia el input de la neurona j en e con la conexión w_{jk}^e)

$$\frac{\partial h_j^e}{\partial w_{jk}^e} = o_k^{e-1}$$

gradiente local

$$\frac{\partial E}{\partial w_{jk}^e} = \frac{\partial E}{\partial h_j^e} o_k^{e-1} = \delta_j^e o_k^{e-1}$$

al mismo tiempo, para los pesos

$$\frac{\partial E}{\partial b_j^e} = \frac{\partial E}{\partial h_j^e} \frac{\partial h_j^e}{\partial b_j^e} \Rightarrow \frac{\partial E}{\partial b_j^e} = \frac{\partial E}{\partial h_j^e} = \delta_j^e$$

E puede ser MSE : $E = \frac{1}{2} \sum_j (y_j^e - o_j^e)^2$ activación

En ese caso:

$$\frac{\partial E}{\partial h_j^e} = \frac{\partial E}{\partial o_j^e} \frac{\partial o_j^e}{\partial h_j^e} \rightarrow o_j^e = g(h_j^e)$$

$$= (y_j^e - o_j^e) g'(h_j^e)$$

$$\Rightarrow \delta_j^e = (y_j^e - o_j^e) g'(h_j^e)$$

$$\Rightarrow \frac{\partial E}{\partial w_{jk}^e} = \underbrace{(y_j^e - o_j^e)}_{\text{comida para la ultima capa}} \cdot o_k^{e-1} g'(h_j^e)$$

$$\Rightarrow \Delta \bar{w} = -\eta \frac{\partial E}{\partial \bar{w}} \quad \Rightarrow \quad \Delta w_{jk}^l = -\eta \frac{\partial E}{\partial w_{jk}^l}$$

$$\Delta w_{jk}^l = -\eta (y_j^l - o_j^l) o_k^{l-1} g'(h_j^l)$$

Para la última neurona capa

$$\frac{\partial E}{\partial h_j^{l-1}} = \sum_m \frac{\partial E}{\partial h_m^l} \frac{\partial h_m^l}{\partial h_j^{l-1}}$$

$$\begin{aligned} h_j^l &= \sum_{k=1}^m w_{jk}^l o_k^{l-1} + b_j^l \\ &= \sum_k w_{jk}^l g(h_k^{l-1}) + b_j^l \end{aligned}$$

$$\frac{\partial h_m^l}{\partial h_j^{l-1}} = w_{mj}^l g'(h_j^{l-1})$$

$$\Rightarrow \frac{\partial E}{\partial h_j^{l-1}} = \sum_m \frac{\partial E}{\partial h_m^l} w_{mj}^l g'(h_j^{l-1}) \quad \frac{\partial E}{\partial h_j^l}$$

$$\begin{aligned} \Delta w_{jk}^{l-1} &= -\eta \frac{\partial E}{\partial w_{jk}^{l-1}} = -\eta \left(\frac{\partial E}{\partial h_j^{l-1}} o_k^{l-2} \right) \\ &= -\eta \left(\sum_m w_{jm}^{l-1} g'(h_j^{l-2}) \frac{\partial E}{\partial h_m^{l-1}} \right) o_k^{l-2} \end{aligned}$$

en notación matricial

última
capa

$$\Delta \underline{w}_{jk}^l =$$

vector
fila

$$\vec{j} \rightarrow \begin{pmatrix} \downarrow k \\ \bar{w} \end{pmatrix}$$

$$- \eta \left(\underline{y}^l - \underline{o}^l \right) g'(\underline{h}^l)$$

PRODUCTO PUNTO A PUNTO (*)
not: vector

$$\begin{pmatrix} j \times 1 \\ \underline{y}, \underline{o} \end{pmatrix} \begin{pmatrix} 1 \times R \\ \underline{o}^T \end{pmatrix}$$

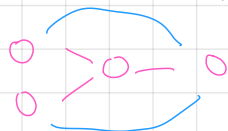
PRODUCTO ESCALAR
(np.dot)

$$\underline{o}^{l-1}$$

fila

$$\Rightarrow \Delta \bar{w} = -\eta (\underline{\bar{y}} - \underline{\bar{o}}) g(\underline{h}) (\underline{o}^{l-1})^T$$

\bar{D}



anteriores
copa

$$\Delta w_{jk}^{l-1} = -\eta \left(\sum_m w_{jm}^l g'(h_j^{l-1}) \frac{\partial E}{\partial h_m^l} \right) o_k^{l-2}$$

$$\sum_m w_{jm}^l \frac{\partial E}{\partial h_m^l} = w_j^T \bar{D} \rightarrow \text{error output}$$

es un prod
matriz x vector
con N^T

$$\Rightarrow \Delta w_{jn}^{l-1} = -\eta w_j^{l-1} \bar{D} g'(h_j^{l-1}) o_k^{l-2}$$

matricial

$$\Rightarrow \Delta w_{jn}^{l-1} = -\eta \left[(w_j^{l-1} \bar{D}) \cdot g'(h) \right] o_k^{l-2}$$

prod pto a pto

$$\Delta \bar{w}^{l-1} = -\eta \left[(w^l \bar{D}) \cdot g'(h) \right] (o^{l-2})^T$$

matricial

$$h_m^l = \sum_s w_{ms} o_s^{l-1} = \sum_s w_{ms} g(h_s^{l-1})$$

$$\rightarrow \frac{\partial h_m^l}{\partial h_k^{l-1}} = w_{mk} g'(h_k^{l-1})$$

