

si η es chico ($\Delta \bar{W} = -\eta \frac{\partial E}{\partial \bar{W}}$) E disminuye

algoritmo de propagación de errores

- inicializar los pesos \bar{W} "ÉPOCAS"
- loop en pasos de aprendizaje (τ fijo o hasta que pare algo)

```

-  $\Delta \bar{W} = 0$ 
- loop entre  $\mu = 1, \dots, P$ 
  -  $\bar{x}_1 \leftarrow \bar{x}^\mu$  (elemento  $x^\mu$  con  $\bar{x}_1$ )
  - loop sobre capas  $m = 1, \dots, M-1$ 
    -  $h_m = \bar{W} \cdot \bar{x}_m$   $\rightarrow$  input de las neuronas en la capa  $m$ 
    -  $x_{m+1} = g(h_m)$   $\rightarrow$  salida de la capa  $m$ 
    -  $\bar{o}^\mu \leftarrow \bar{x}_M$ 
    -  $\delta_m = -\frac{\partial E}{\partial \bar{o}^\mu} g'(h_m)$ 
    - loop sobre capas ( $m = m, \dots, 2 \rightarrow$  para atrás)
      -  $\delta_{m-1} = [\bar{W}^T \delta_m] \otimes g'(h_{m-1})$ 
    - loop sobre capas  $m = 1, \dots, m-1$ 
      - (*)  $\Delta \bar{W} += \eta \delta_{m+1} \bar{x}_m$  poco complicado  $\hookrightarrow$  hay que guardar los  $x$ 
  -  $\bar{W} = \bar{W} + \Delta \bar{W}$ 

```

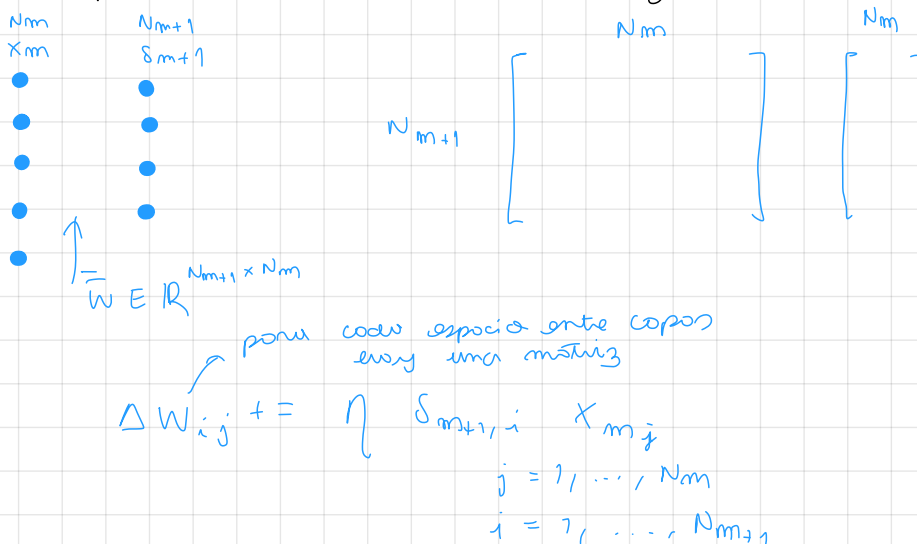
hay que guardar x, h, δ, W y ΔW

CONTROL DE DETERMINACIÓN

$\Delta W < w_0 \rightarrow$ break

¿Cuál es el objetivo de w_0 ?

Por componentes está en el Hertz



uno puede alterar los umbrales poniendo una neurona adicional

$$x_{m,i} = g \left(\sum_{j=1}^m w_{ij} x_{m,j} - \theta_i \right)$$

podemos poner una neurona con una salida constante, que nos de estos parámetros

↓
esta neurona estaría en la capa de entrada y le mandaría señales a todos)

no hay un teorema que nos garantice la convergencia

Heurística

inicialización de los pesos

si los pesos son muy grandes

↓
derivada pequeña, g' es lo que ΔW pequeño ($|\Delta W| < 1$)

↓
modificaciones muy pequeñas

- se quiere que $|\bar{w} \cdot \bar{x}_0| \sim \underbrace{O(1)}_{\text{orden 1}}$

↓
normalizar los inputs

(mucha de pto, se puede poner con media 0 y varianza $1/N \rightarrow$ no se pierde info)

↓
esto por defecto en varios paquetes)

N : dim de la capa de entrada

h va a tener varianza

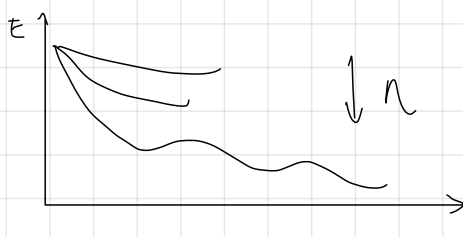
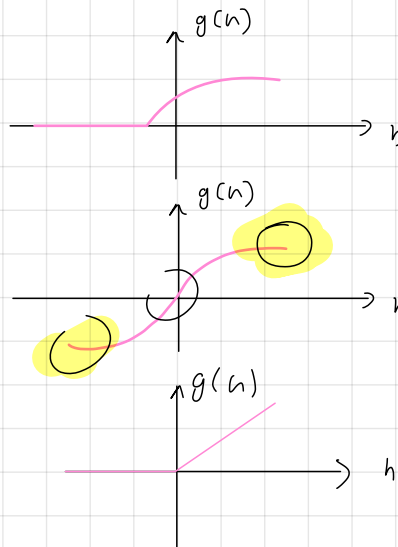


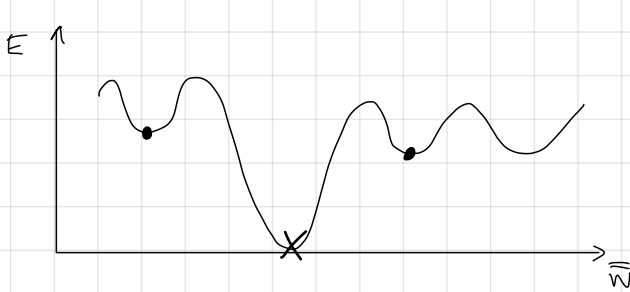
grafico lo que se quiere minimizar para los distintos valores de N

tende a fallar por ser un método local

↓
el método puede quedar atrapado en mínimos locales o pto silla no deseado)

dos típicos
sigmoide





Para ver si nos quedamos en un min local:

- se pueden probar muchos CF distintos
- o bien modifican el método de aprendizaje

Podemos pensar que se trata de una pelota \Rightarrow si tiene impulso atravesará los mínimos locales

$E(\bar{w})$

$$\bar{g} = \nabla_{\bar{w}} E(\bar{w}) \rightarrow \text{MOMENTO}$$

$$\bar{v}_t = \beta \bar{v}_{t-1} + (1-\beta) \bar{g}_t$$

\downarrow
promedio del gradiente con una tasa de tiempo de $1/\beta$

o

o

\leftarrow

$$\Delta w = -\eta \bar{v}_t$$

guarda en memoria la dirección en la que es el movimiento

otra forma

NESTERON

$$\bar{v}_{t+1} = \beta \bar{v}_{t-1} + \eta \nabla E(\bar{w} - \beta \bar{v}_{t-1})$$

$$\Delta \bar{w} = -\bar{v}_t$$

ADAM

$$\bar{m}_t = \beta \bar{m}_{t-1} + (1-\beta) \bar{g}_t$$

$$\bar{v}_t = \beta \bar{v}_t + (1-\beta) |\bar{g}_t|^2$$

$$\Delta \bar{w} = -\eta \frac{\bar{m}_t}{\sqrt{\bar{v}_t} + \epsilon}$$

stochastic gradient descent

SGD (porá metros.

nº de elementos de un Batch)

$\Delta \bar{w} \rightarrow$ suma sobre $\mathcal{H} = 1, \dots, p$

Podemos dividir el conjunto en subconjuntos



$$\Delta w = \sum_{\mathcal{H}=1}^p \Delta \bar{w}^{\mathcal{H}}$$

$$= \sum_{b=1}^B \left[\sum_{\mathcal{H}_b=1}^{p/B} \Delta \bar{w}^{\mathcal{H}_b} \right] \frac{1}{B}$$

cada uno de estos valores es una variable aleatoria que alterna entre el vector Δw