



# **REDES NEURONALES 2024**

**Clase 17 parte 1**

**Jueves 17 de octubre 2024**

**FAMAF, UNIVERSIDAD NACIONAL DE CÓRDOBA**

**INSTITUTO DE FÍSICA ENRIQUE GAVIOLA (UNC-CONICET)**

M24/10/23 c18p2

## El método del descenso por el gradiente

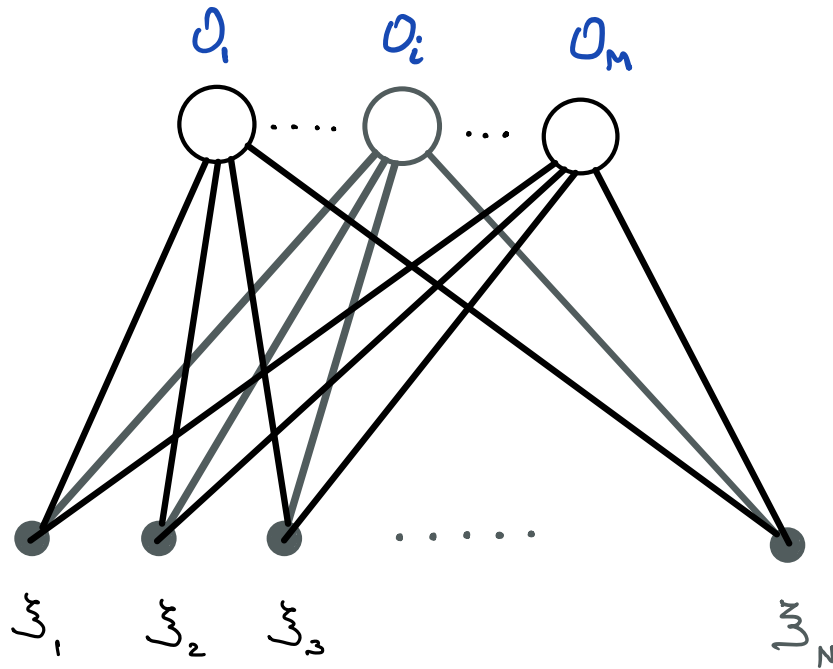
Hoy presentaremos por primera vez en este curso un elemento central en el aprendizaje supervisado moderno, el cual nos acompaña todos los días en nuestras investigaciones y nuestras aplicaciones de IA. Veremos en esta parte la estrecha relación existente entre aprender y optimizar una función, o sea, encontrar sus valores mínimos, locales y globales.

Definir como función error o función costo (loss function, en inglés).

$$E(\bar{w}) = \sum_{\mu=1}^P \left( \sum_{i=1}^M \frac{1}{2} (Y_i^{\mu} - O_i^{\mu})^2 \right)$$

Suma sobre la p  
elementos del conjunto  
de entrenamiento

Suma sobre la M  
neuronas de la capa de  
salida



Cuando la presentamos la entrada  $\mu$  obtenemos el resultado dado por la regla del perceptrón lineal:

$$O_i^{\mu} = g(h_i^{\mu}) = h_i^{\mu} = \sum_k^N w_{ik} \xi_k^{\mu} = \underbrace{\bar{w}_i \cdot \bar{\xi}^{\mu}}_1$$

La idea es comenzar con un vector inicial arbitrario  $W$ , como hacíamos en el caso de la neurona de salida binaria y presentarle un ejemplo, el  $\mu$  del conjunto de entrenamiento:

$$\overrightarrow{\nabla E} : \mathbb{R}^{N \times M} \xrightarrow{\mathbb{R}^{N \times M}} \mathbb{R}$$

$$\nabla E_{ik} = \frac{\partial E}{\partial w_{ik}}$$

$$w_{ik}^{\text{nuevo}} = w_{ik}^{\text{anterior}} + \Delta w_{ik}$$

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

$$= -\eta \frac{\partial}{\partial w_{ik}} \frac{1}{2} \sum_{\mu=1}^P \left( y_i^{\mu} - o_i^{\mu} \right)^2$$

$$= -\eta \sum_{\mu=1}^P \frac{1}{2} \frac{\partial}{\partial w_{ik}} \left( y_i^{\mu} - \sum_j w_{ij} z_j^{\mu} \right)^2$$

$$= -\eta \sum_{\mu=1}^P \frac{2}{2} \left( y_i^{\mu} - \sum_j w_{ij} z_j^{\mu} \right) \frac{\partial}{\partial w_{ik}} \left( - \sum_j w_{ij} z_j^{\mu} \right)$$

$$= \eta \sum_{\mu=1}^P \left( y_i^{\mu} - o_i^{\mu} \right) z_k^{\mu}$$

$$\Delta w_{ik} = \eta \sum_{\mu=1}^P \left( y_i^{\mu} - o_i^{\mu} \right) z_k^{\mu}$$

$$w_{ik}^{\text{nuevo}} = w_{ik}^{\text{anterior}} + \Delta w_{ik}$$

Si definimos:

$$\delta_i^H = \sum_i^H - O_i^H$$

entonces:

$$W_{ik}^{\text{nuevo}} = W_{ik}^{\text{anterior}} + \eta \sum_{H=1}^P \delta_i^H \sum_k^H$$

Esta regla nos permite buscar un mínimo local de una función de muchas variables. En nuestro caso buscamos el mínimo de la función *Error Cuadrático Medio* sobre todos los posibles valores de las sinapsis y umbrales. En otras palabras, se trata de encontrar el juego de sinapsis y umbrales que produce el menor de aprendizaje sobre el conjunto de entrenamiento.

Esta regla fue deducida y bautizada en forma independiente varias veces:

- Regla ADELIN o Delta (Widrow y Hoff, 1960)
- Regla Least Means Square (1986)
- Regla Rescola wagner (1976)

Algoritmo de actualización en línea (pseudo código)

Elegimos lo  $N$  parámetros  $W$  al azar

Sea  $\text{época}=1$

Repetimos hasta que  $\text{época}=\text{época\_maxima}$

Sea  $\mu = 1$

Repetimos, hasta que  $\mu = p$

$O_i^h = \bar{w}_i \cdot \tilde{z}^h$   $i = 1, 2, \dots, M$   
Calculamos la loss  $E$  para el ejemplo  $\mu$  y  
obtenemos los  $M$  valores de salida  $O$

$\Delta w_{ik} = -\eta \frac{\partial E}{\partial w_{ik}} = \eta \delta_i^h \tilde{z}_k^h$   
Calculamos los  $M \times N$  valores de los incrementos  
de cada parámetro

$$\bar{w}_i = \bar{w}_i + \Delta \bar{w}_i$$

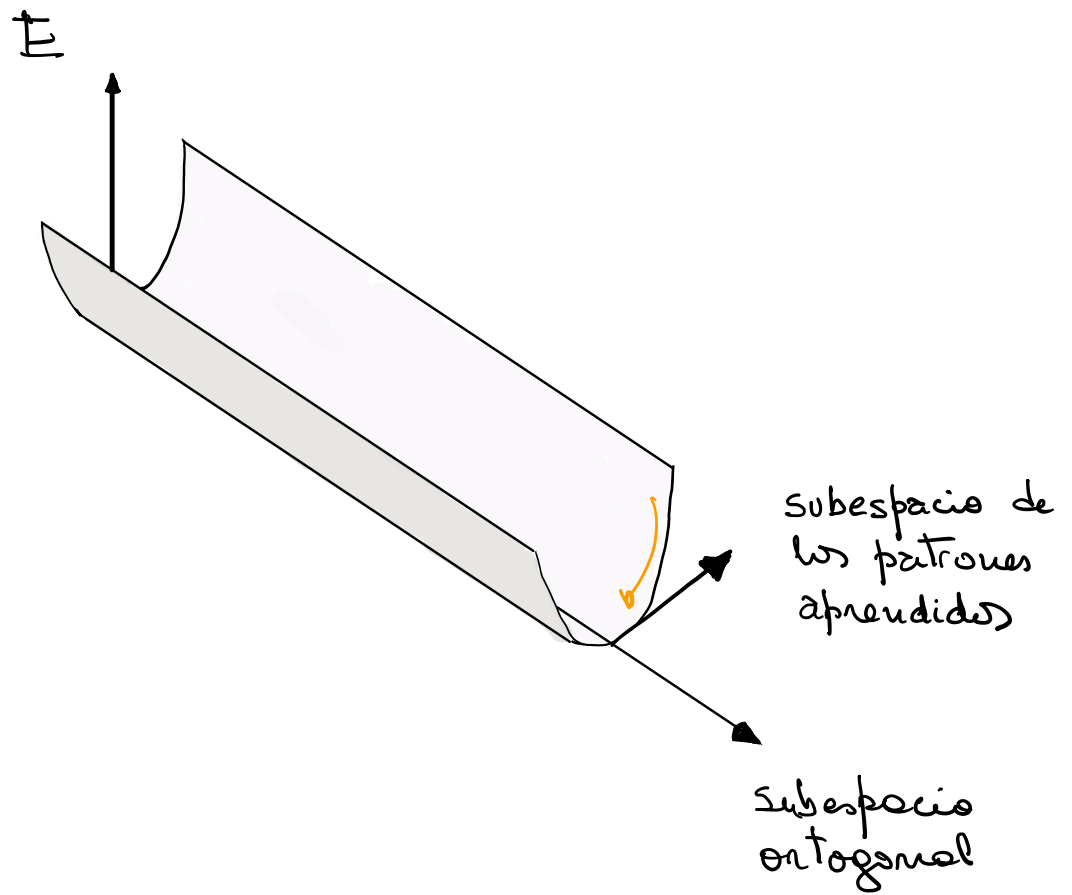
Actualizamos los  $M \times N$

$\mu = \mu + 1$

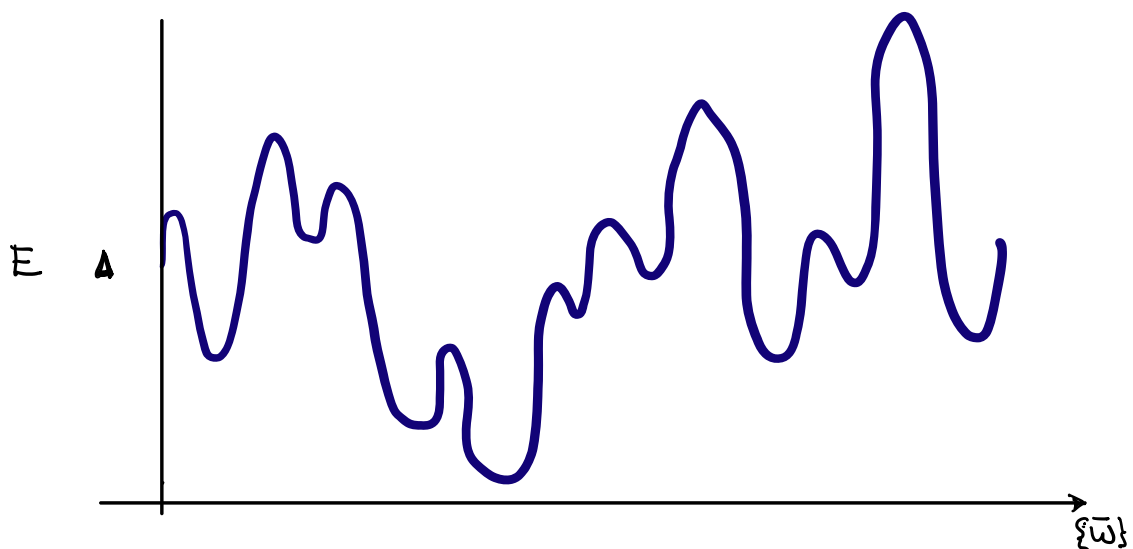
Volvemos a “Repetimos ( $\mu$ )”

$\text{época} = \text{época} + 1$

Volvemos a “Repetimos ( $\text{época}$ )”

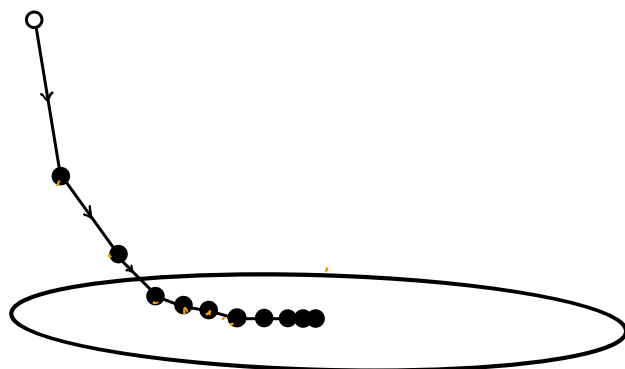


Ahora no vale la separabilidad lineal. Tenemos que intentar aprender descendiendo por la compleja función loss, la cual es bastante

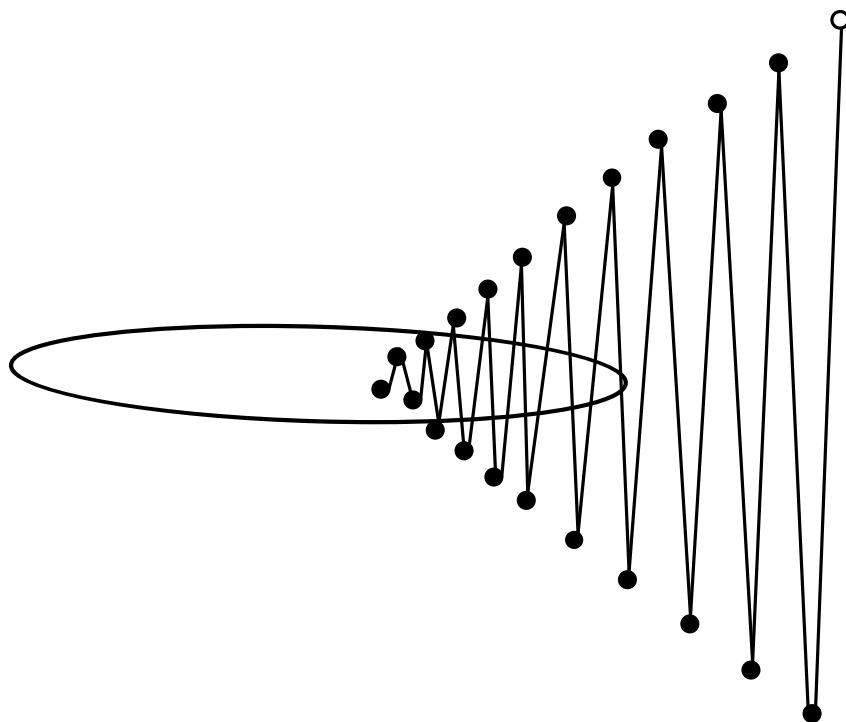


Ejemplo:

$$E = x^2 + 20y^2$$



$$\eta = 0.02$$



$$\eta = 0.476$$



