



# **REDES NEURONALES 2024**

**Clase 18 parte 1**

**Jueves 23 de octubre 2024**

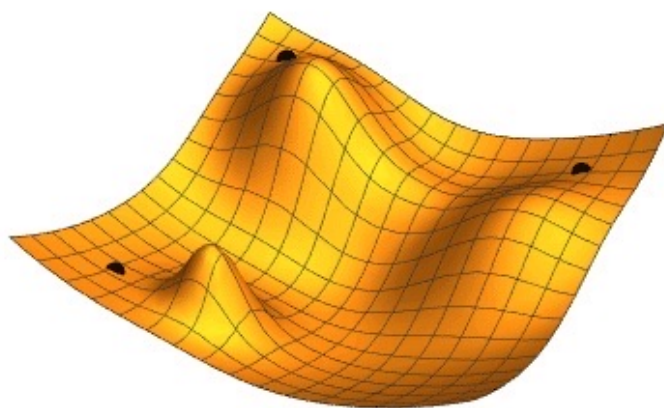
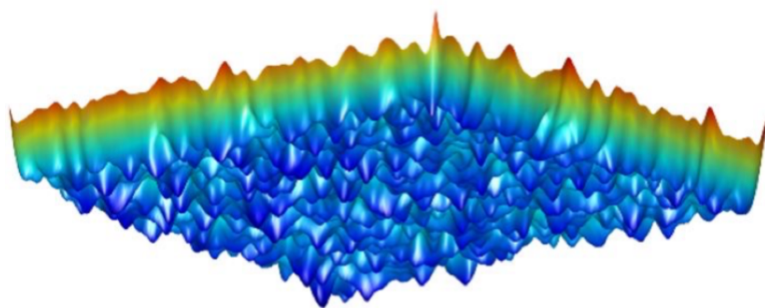
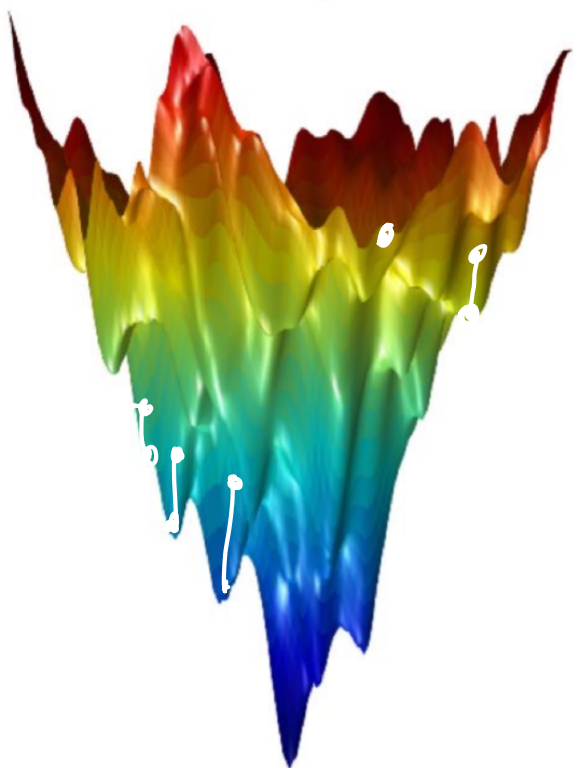
**FAMAF, UNIVERSIDAD NACIONAL DE CÓRDOBA**

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

J26/10/23 c19

Aunque quizá no sea obvio, debemos saber que el Método de Descenso por el Gradiente es muy limitado como método de optimización global de la función loss  $E(\{W\})$ . Veamos porqué.

- 1.- La función loss o error tiene ahora muchísimos mínimos locales debido a la no-linealidad de las funciones de activación sumado al hecho de que se elevan al cuadrado y a la alternancia y aleatoriedad de los signos de las sinapsis.
- 2.- El método del descenso por el gradiente (MDG) se queda atrapado en el primer mínimo local que encuentra, y es muy poco probable que este sea el que tiene el menor valor de la función loss  $E$ .
- 3.- Las funciones de activación exponenciales y sus derivadas son lentas de calcular en punto flotante.
- 4.- El MDG depende fuertemente del valor inicial aleatorio de los parámetros del problema (las sinapsis y el umbral) debido a la rugosidad de  $E(\{W\})$ .
- 5.- El MDG es muy sensitivo al valor de la razón de aprendizaje y este no puede determinarse en forma óptima a priori.
- 6.- El MDG trata de la misma forma a todas las direcciones en el espacio  $N$  dimensional en el cual viven los vectores  $\overline{W}$ . Nos gustaría tener un método que dé pasos cortos en direcciones íngrimes del paisaje rugoso de la función  $E(\{\overline{W}\})$ .
- 7.- El MDG requiere de tiempos exponenciales en  $N$  para escapar de los puntos de ensilladura, los cuales proliferan a medida que  $N$  crece.



En realidad podemos buscar un  $\eta$  óptimo  $\eta_{\text{opt}}$ , pero eso es muy caro en cálculo computacional. Para ello necesitamos realizar una expansión de Taylor de la función  $E(\bar{W})$  alrededor del mínimo. Nos paramos en un punto cualquiera de  $\mathbb{R}^{n \times n}$  que denotamos  $\bar{W}$ .

$$E(\bar{w} + \bar{r}) = E(\bar{w}) + \partial_{\bar{w}} E(\bar{w}) \bar{r} + \frac{1}{2} \partial_{\bar{w}}^2 E(\bar{w}) \bar{r}$$

Sea  $W_{\min}$  el punto que minimiza  $E(W)$ .

$$\bar{W} = \bar{W}_{\min} + \bar{r}$$

Si es un mínimo, todas las componentes del gradiente son nulas:

$$\partial_{\bar{w}} E(\bar{w}) = 0$$

Derivando con respecto a  $\bar{r}$

$$\frac{\partial E}{\partial \bar{w}} \cdot \frac{\partial(\bar{w} + \bar{r})}{\partial \bar{r}} = \frac{\partial E}{\partial \bar{w}} + \frac{\partial E}{\partial \bar{w}} + 2 \left( \frac{\partial^2 E}{\partial \bar{w}^2} \right) \bar{r}$$

$$\bar{v} = - \left( \frac{\partial E}{\partial \bar{w}} \right) \frac{1}{\left( \frac{\partial^2 E}{\partial \bar{w}^2} \right)}$$

$$\bar{w} = \bar{w}_{\min} + \bar{v}$$

$$= \bar{w}_{\min} + \frac{1}{\left( \frac{\partial^2 E}{\partial \bar{w}^2} \right)} \left( \frac{\partial E}{\partial \bar{w}} \right)$$

$$\eta_{\text{opt}} = \frac{1}{\left( \frac{\partial^2 E}{\partial \bar{w}^2} \right)}$$

**ESTO ES CARÍSIMO EN TÉRMINOS COMPUTACIONALES**

Si fuésemos muy formales

$$E(\bar{\omega} + \bar{v}) \approx E(\bar{\omega}) + \nabla_{\bar{\omega}} E(\bar{\omega}) \cdot \bar{v} + \frac{1}{2} \bar{v}^T H(\bar{\omega}) \bar{v}$$

donde  $H$  es el Hessiano de  $E$

Si estamos en un mínimo  $\nabla E(\bar{\omega} + \bar{v})$

Derivando con respecto a  $\bar{v}$

$$\nabla_{\bar{\omega}} E(\bar{\omega}) = -H(\bar{\omega}) \bar{v}_{\text{opt}}$$

$$\bar{v}_{\text{opt}} \approx \bar{v}_t = H^{-1}(\bar{\omega}_t) \nabla_{\bar{\omega}} E(\bar{\omega}_t)$$

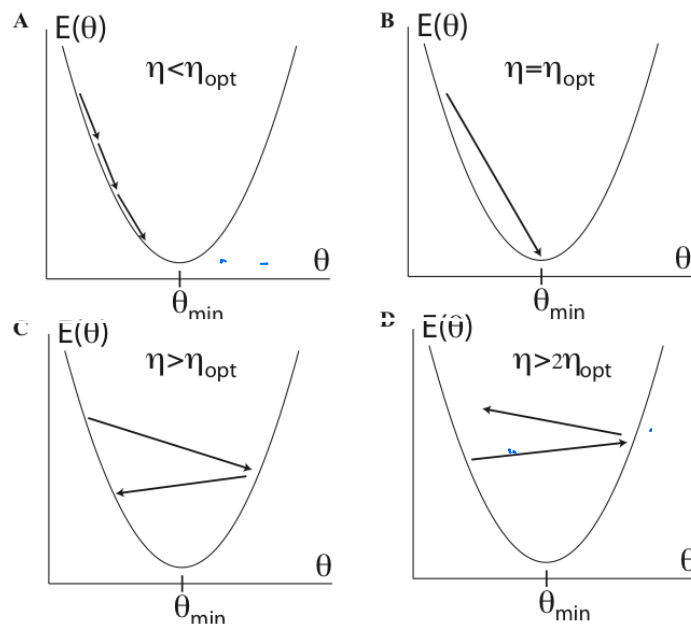
Si estamos en un mínimo de  $E$  :  $\bar{\nabla}_{\omega} E(\bar{\omega} + \tilde{v}) = 0$

Derivando con respecto a  $\tilde{v}$

$$\nabla_{\bar{\omega}} E(\bar{\omega}) = - H(\bar{\omega}) \bar{v}_{\text{opt}}$$

$$\bar{v}_{\text{opt}} \approx \tilde{v}_t = H^{-1}(\bar{\omega}_t) \nabla_{\omega} (E(\bar{\omega}_t))$$

$$\eta < \frac{2}{\lambda_{\max}}$$



**FIG. 8 Effect of learning rate on convergence.** For a one dimensional quadratic potential, one can show that there exists four different qualitative behaviors for gradient descent (GD) as a function of the learning rate  $\eta$  depending on the relationship between  $\eta$  and  $\eta_{\text{opt}} = [\partial_{\theta}^2 E(\theta)]^{-1}$ . (a) For  $\eta < \eta_{\text{opt}}$ , GD converges to the minimum. (b) For  $\eta = \eta_{\text{opt}}$ , GD converges in a single step. (c) For  $\eta_{\text{opt}} < \eta < 2\eta_{\text{opt}}$ , GD oscillates around the minima and eventually converges. (d) For  $\eta > 2\eta_{\text{opt}}$ , GD moves away from the minima. This figure is adapted from (LeCun *et al.*, 1998b).

