

UNIDADES NO LINEALES

Consideremos el caso en el cual la neurona de salida tiene una función de activación NO LINEAL

$g(z)$.

Podemos aplicar la misma regla de descenso por el gradiente, haciendo pequeños cambios en los pesos sinápticos, en la dirección contraria al gradiente de la función error.

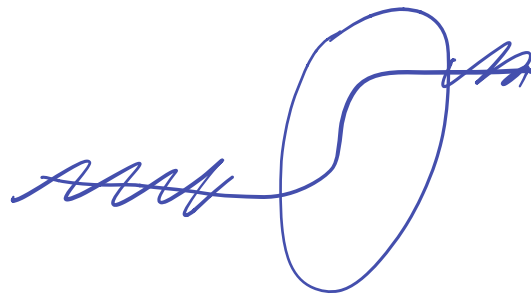
$$\begin{aligned}
 E(\bar{w}) &= \frac{1}{2} \sum_i^M \sum_{\mu}^P (y_i^{\mu} - o_i^{\mu})^2 \\
 &= \frac{1}{2} \sum_i \sum_{\mu} (y_i^{\mu} - g(h_i^{\mu}))^2 \\
 &= \frac{1}{2} \sum_i \sum_{\mu} \left(y_i^{\mu} - g\left(\sum_k w_{ik} z_k^{\mu}\right) \right)^2
 \end{aligned}$$

$$\frac{\partial E}{\partial w_{ik}} = \frac{1}{2} \sum_i \sum_{\mu} 2 (y_i^{\mu} - g(h_i^{\mu})) \left(-g'(h_i^{\mu}) \right) \frac{dh_i^{\mu}}{dw_{ik}}$$

$$\frac{\partial E}{\partial w_{ik}} = \frac{1}{2} \sum_i \sum_{\mu} 2 (y_i^{\mu} - g(h_i^{\mu})) (-g'(h_i^{\mu})) \sum_k^{\mu}$$

$$\frac{\partial E}{\partial w_{ik}} = - \sum_i \sum_{\mu} (y_i^{\mu} - 0_i^{\mu}) g'(h_i^{\mu}) \sum_k^{\mu}$$

$$\frac{\partial E}{\partial w_{ik}} = - \sum_{\mu} \delta_i^{\mu} \sum_k^{\mu}$$

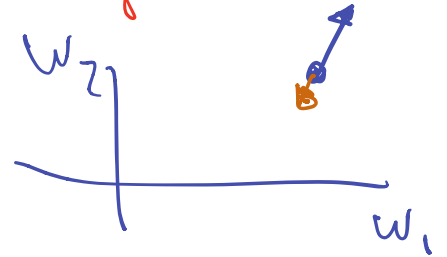


$$\delta_i^{\mu} \equiv [y_i^{\mu} - 0_i^{\mu}] g'(h_i^{\mu})$$

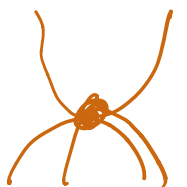
$$\Delta w_{ik} = \eta \delta_i^{\mu} \sum_k^{\mu}$$

Ahora podemos hacer lo mismo que con las neuronas de salida lineales: ENTRENAR LA RED COMENZANDO CON SINAPSIS $\{\bar{w}_i\}$ ALEATORIAS

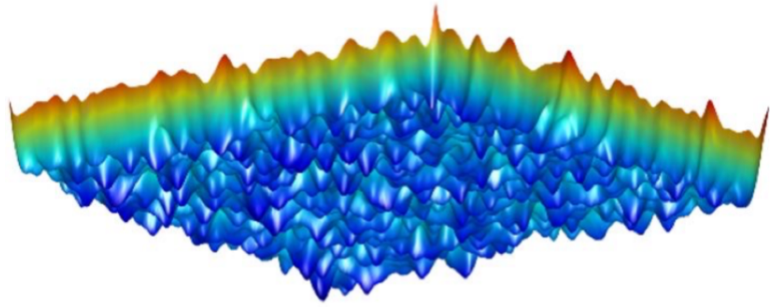
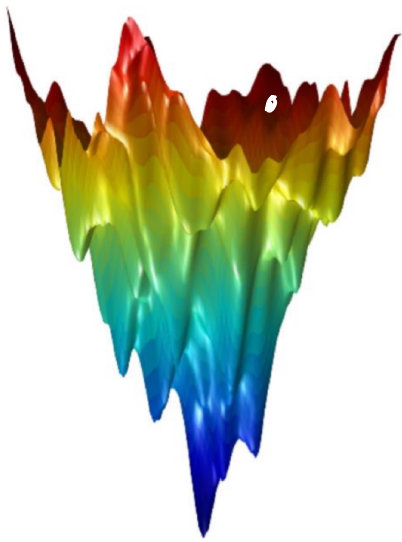
Los problemas del descenso por el gradiente



- La función error (o costo) tiene muchísimos mínimos locales, debido a la forma cuadrática de funciones fuertemente no lineales, a la alternancia de signos de los w_{ik} y su aleatoriedad.
- El método de descenso por el gradiente (MDG) se puede atropado en mínimos locales con valores altos de E .
- Las funciones no lineales son lentas de evaluar, y más lento es evaluar sus derivadas $g'(h)$.
- El MDG es muy sensitivo al valor de la razón de aprendizaje η , y no es posible conocer exactamente su valor óptimo.



- El MDG depende FUERTEMENTE del valor inicial de los parámetros $\{\bar{w}_i\}$, por culpa de la rugosidad de E
- El MDG trata las $M \times N$ direcciones en \mathbb{R}^N de la misma forma (con el mismo η), siendo que no todas las direcciones son igualmente importantes. Nos gustaría dar "pasos" largos en las direcciones en las cuales la función E varía poco, y "pasos" cortos en las que varía mucho.
- El método de descenso por el gradiente requiere de tiempos exponencialmente largos para poder escapar de puntos de ensilladura.



Podemos encontrar un γ_{opt} . Nos fijamos en un punto cualquiera de \mathbb{R}^{N+M} y denotemos por $\bar{w} \in \mathbb{R}^{N+M}$ a ese punto

$$E(\bar{w} + \bar{v}) = E(\bar{w}) + \underbrace{\partial_{\bar{w}} E(\bar{w})}_{\bar{w}'} \bar{v} + \frac{1}{2} \partial_{\bar{w}}^2 E(\bar{w}) \bar{v}$$

Sea \bar{w}_{min} el valor de \bar{w} que minimice

$E(\bar{w})$.

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

Por ser mínimo

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

Derivando con respecto a \bar{v}

$$\frac{\partial E}{\partial \bar{\omega}} \cdot \frac{\partial \bar{\omega} + \vec{v}}{\partial v} = \frac{\partial E}{\partial \bar{\omega}} + \frac{\partial E}{\partial \bar{\omega}} + 2 \left(\frac{\partial^2 E}{\partial \bar{\omega}^2} \right) \vec{v}$$

$$\vec{v} = - \left(\frac{\partial E}{\partial \bar{\omega}_m} \right) \cdot \frac{1}{\left(\frac{\partial^2 E}{\partial \bar{\omega}^2} \right)}$$

$$\bar{\omega} = \bar{\omega}_{min} + \vec{v}$$

$$= \bar{\omega}_{min} + \left[\frac{1}{\left(\frac{\partial^2 E}{\partial \bar{\omega}^2} \right)} \right] \cdot \frac{\partial E}{\partial \bar{\omega}}$$

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

oo Es es CRÍSIMO !!
||

Si fuéramos muy formales

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

H : Hessiano de E

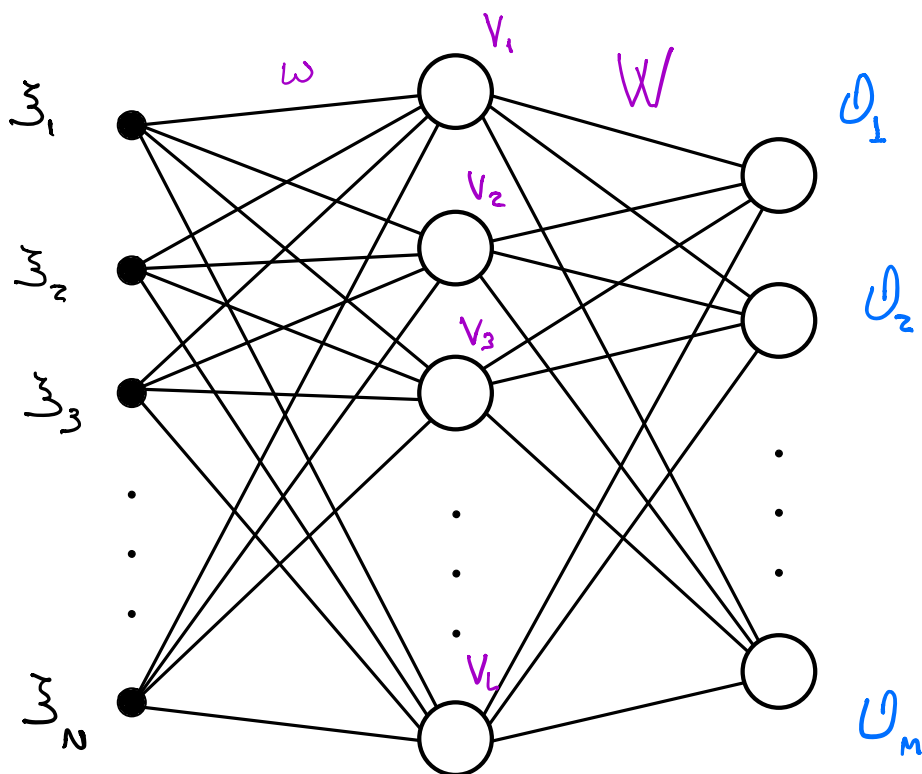
Si estamos en un mínimo $\nabla_{\bar{w}} E(\bar{w} + \bar{v}) = 0$

Derivando con respecto a \bar{v}

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

$$\bar{v}_{opt} \approx \bar{v}_t = -H^{-1}(\bar{w}_t) \nabla_w (E(\bar{w}_t))$$

BACK PROPAGATION PARA UNA CAPA OCULTA



$$h_j^H = \sum_{k=1}^N w_{jk} z_k^H$$

$$v_j^H = g_j(h_j^H) = g_j\left(\sum_{k=1}^N w_{jk} z_k^H\right)$$

$$h_i^\mu = \sum_{j=1}^L W_{ij} V_j$$

$$O_i^\mu = g_2(h_i^\mu) = g_2\left(\sum_j W_{ij} V_j^\mu\right)$$

$$= g_2\left(\sum_j W_{ij} \cdot g_1\left(\sum_{k=1}^M W_{jk} \sum_k^\mu\right)\right)$$

$$E[\{\bar{w}\}] = \frac{1}{2} \sum_{\mu=1}^P \sum_{i=1}^M [y_i^\mu - O_i^\mu]^2$$

$$= \frac{1}{2} \sum_{\mu=1}^P \sum_{i=1}^M \left[y_i^\mu - g_2\left(\sum_j W_{ij} \cdot g_1\left(\sum_{k=1}^M W_{jk} \sum_k^\mu\right)\right) \right]^2$$

$$W_{ij}^{\text{nuevo}} = W_{ij}^{\text{anterior}} + \Delta W_{ij}$$

$$\begin{aligned}
 \Delta W_{ij} &= -\eta \frac{\partial E}{\partial W_{ij}} \\
 &= \eta \sum_{\mu=1}^P \cancel{\frac{\partial}{\partial}} [y_i^\mu - O_i^\mu] g'(h_i) V_j^\mu \\
 &= \eta \sum_{\mu=1}^P \delta_i^\mu V_j^\mu
 \end{aligned}$$

$$\delta_i^\mu = g'(h_i) [y_i^\mu - O_i^\mu]$$

Entonces, presentamos el elemento μ del conjunto de entrenamiento, \vec{x}^μ , al perceptron con una capa oculta, y obtenemos el resultado \vec{O}^μ .

Con \vec{O}^μ calculamos el error cuadrático $E(\{\vec{O}\})$,

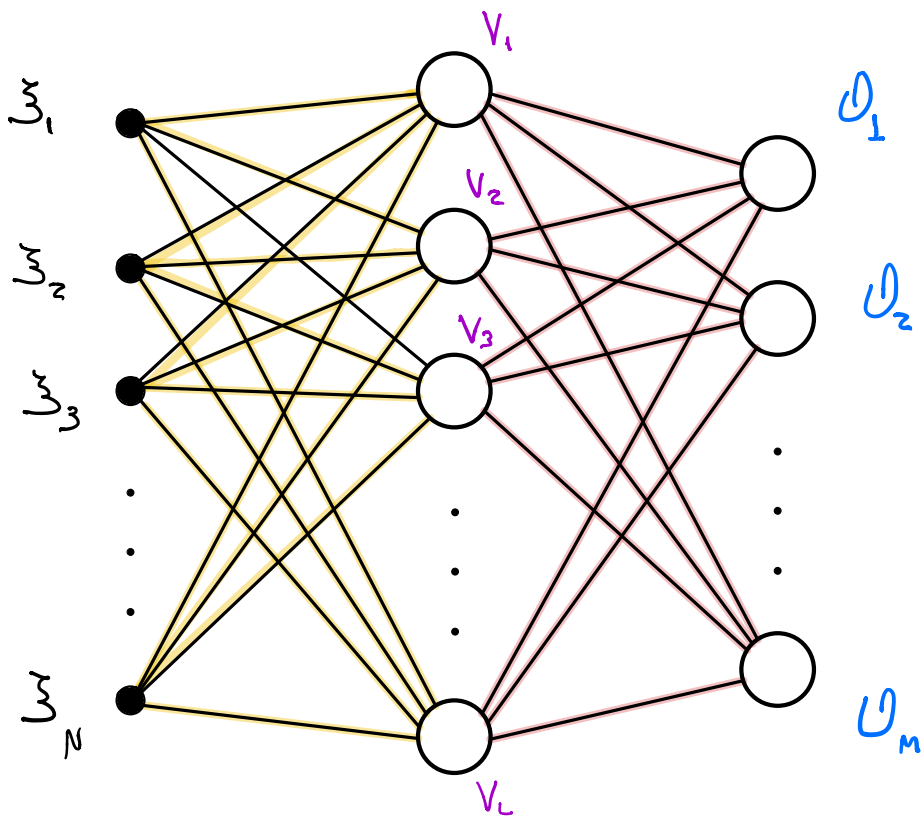
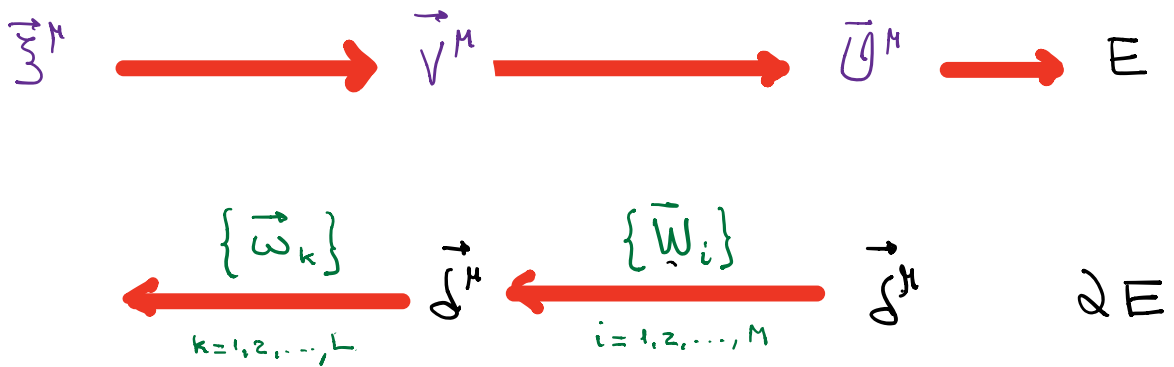
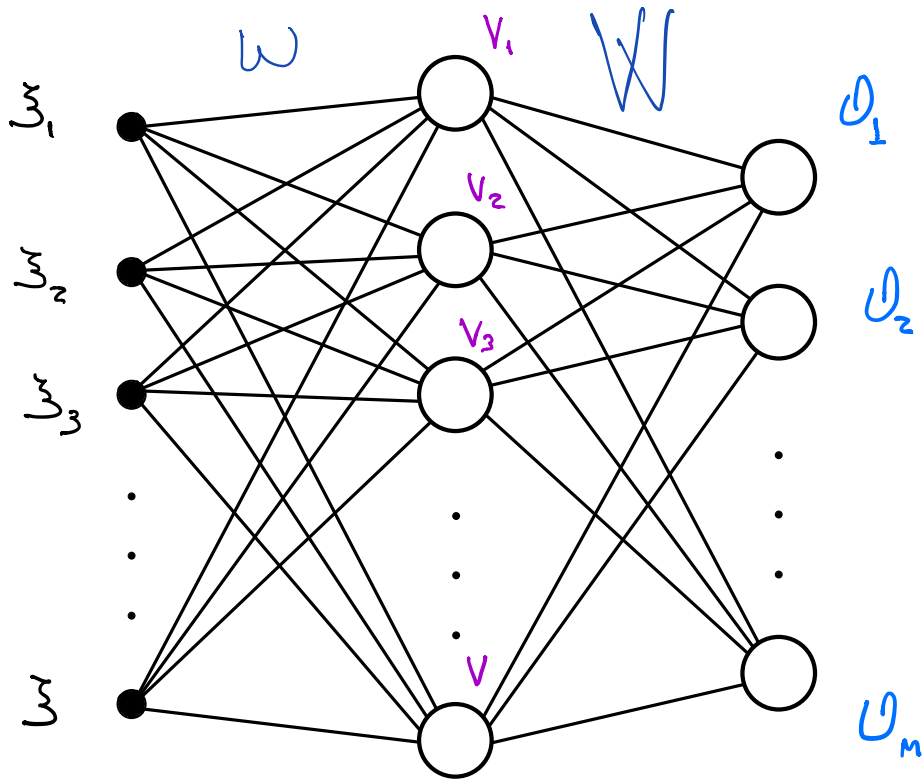
$$\Delta \omega_{jk} = - \eta \frac{\partial E}{\partial \omega_{jk}}$$

$$= \eta \sum_{\mu}^P [y_i^{\mu} - o_i^{\mu}] g'_1(h_i^{\mu}) w_{ij} g'_2(h_j^{\mu}) z_k^{\mu}$$

$$= \eta \sum_{\mu}^P \delta_i^{\mu} w_{ij} g'_2(h_j^{\mu}) z_k^{\mu}$$

$$= \eta \sum_{\mu}^P \delta_j^{\mu} z_k^{\mu}$$

$$\delta_j^{\mu} = g'_2(h_j^{\mu}) \sum_i w_{ij} \delta_i^{\mu}$$



¿cuantos suplementos tenemos?

$$(N \times L) + (L \times M) = L \times (N + M)$$

EL ALGORITMO

- presento \vec{z}^M , calculo \vec{v}^M , luego \vec{O}^M .
- con \vec{O}^M y \vec{z}^M calculamos $E(\{\bar{w}\})$
- con E calculo los δ_i^M ($i=1, \dots, M$) y con ellos actualizo los pesos W_{ij}

$$W_{ij}^{\text{nuevo}} = W_{ij}^{\text{anterior}} + \Delta W_{ij}$$

con los δ_i^M calculo los δ_j^M ($j=1, 2, \dots, L$) y con ellos actualizemos los pesos ω_{jk}

$$\omega_{jk}^{\text{nuevo}} = \omega_{jk}^{\text{anterior}} + \Delta \omega_{jk}$$

EL ALGORITMO

- $\mu = 1$

WHILE ($\mu \leq P$) repeat

1.- presento \vec{z}^{μ} , calculo \vec{v}^{μ} , luego \vec{O}^{μ} .

2.- con \vec{O}^{μ} y \vec{z}^{μ} calculamos $E(\{\bar{w}\})$

3.- con E calculo los δ_i^{μ} ($i=1, \dots, M$) y con ellos actualizo los pesos W_{ij}

$$W_{ij}^{\text{nuevo}} = W_{ij}^{\text{anterior}} + \Delta W_{ij}$$

$$\Delta W_{ij} = \sum W_{ij}^M$$

4.- con los δ_i^{μ} calculo los δ_j^{μ} ($j=1, 2, \dots, L$) y con ellos actualizemos los pesos w_{jk}

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

$\mu = \mu + 1$

Hay diferentes formas de hacer el proceso

- **en línea**: luego de presentar el ejemplo redefinimos todos los acoplamientos sinápticos.
- **en batch**: luego de presentar el ejemplo almacenamos la variación de los acoplamientos sinápticos pero actualizo cuando termine de mostrarle los P elementos del conjunto de entrenamiento.

$$\Delta W_{ij} = \Delta W_{ij}^1 + \Delta W_{ij}^2 + \dots + \Delta W_{ij}^P$$

$$\Delta \omega_{jk} = \Delta \omega_{jk}^1 + \Delta \omega_{jk}^2 + \dots + \Delta \omega_{jk}^P$$

Veremos que hoy se usa la técnica MINIBATCHES

En la próxima clase veremos MUCHOS trucos para superar, heurísticamente, la dificultad intrínseca de tener que minimizar una función intrínseca y compleja como es E .

