

REDES NEURONALES CON PYTHON

Francisco Tamarit

El algoritmo de retro propagación del error

Recordemos que una única neurona tiene una gran capacidad de modelar.

- Perceptrón simple con función de activación binaria, apto para clasificar en dos categorías (presentamos el primer algoritmo de aprendizaje automático).
- Perceptrón simple con función de activación lineal, apto para hacer **regresión lineal** y resolver problemas linealmente independientes (muy muy simples).
- Perceptrón simple con función de activación no lineal, apto para hacer una **regresión logística** y aprender a resolver problemas un poco más complicados, pero no mucho más complicadas.

Recordemos que estamos analizando el caso de **aprendizaje supervisado**, para lo cual necesitamos un conjunto de entrenamiento C que resume la experiencia previa:

$$C = \left\{ \left(\begin{pmatrix} \bar{x}^1 \\ \vdots \\ \bar{x}^N \end{pmatrix}, \begin{pmatrix} \bar{y}^1 \\ \vdots \\ \bar{y}^N \end{pmatrix} \right) \right\}$$

, P

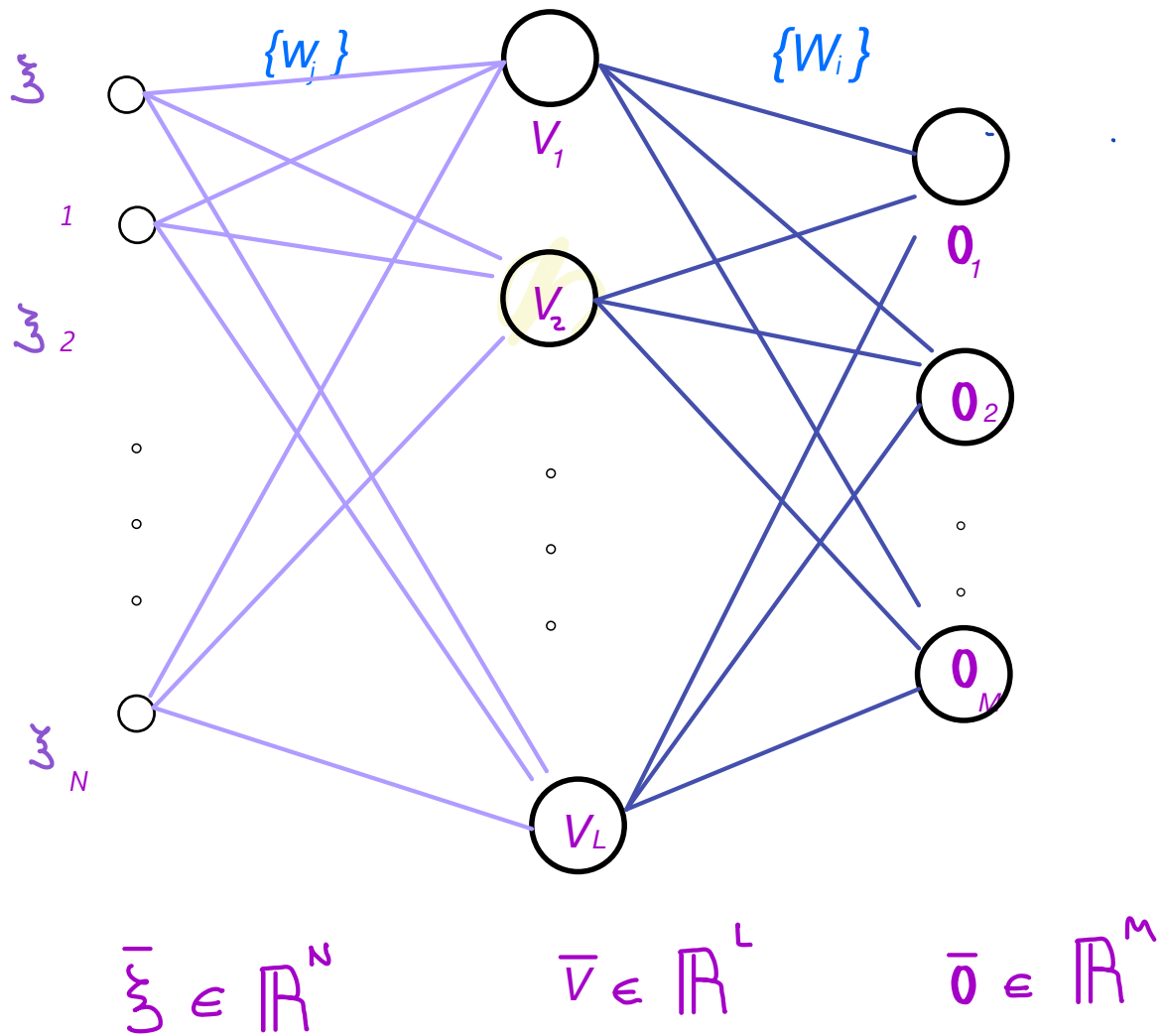
Entradas

Salidas etiquetadas

Nos limitaremos por ahora al caso de regresiones, o sea, de salidas reales. Todas las funciones de activación serán no lineales.

Nos limitaremos también a una arquitectura de red neuronal muy simple y eficiente, que permite tratar problemas sofisticados con mucho éxito. Se trata de las famosas *Redes Neuronales Feed-Forward*, o, en español, redes neuronales alimentadas hacia adelante. .

Por razones didácticas, antes de analizar el caso más general de un número arbitrario de capas de neuronas, analizaremos en detalle el caso de una *red feed-forward con una única capa oculta*.



Queremos que la red aprenda una función objetivo f que no conocemos y cuya expresión analítica no pretendemos conocer.

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

Este es un diagrama esquemático de una red neuronal feed-forward con N entradas, M salidas y una única capa oculta de L neuronas.

- Usaremos la letra griega ξ_k para indicar la k -ésima entrada y reservamos la letra k para nombrarlas.
- Llamaremos V_j a la función de activación de las neuronas de la capa oculta y usaremos la letra j para nombrarlas.
- Llamaremos O_i a las neuronas de la capa de salida y usaremos la letra i para nombrarlas.

Pensamos en términos de vectores:

$$\bar{\xi} \in \mathbb{R}^N$$

$$\bar{V} \in \mathbb{R}^L$$

$$\bar{O} \in \mathbb{R}^M$$

$$\bar{W}_i \in \mathbb{R}^L$$

$$\bar{W}_j \in \mathbb{R}^N$$

El proceso de entrenamiento es similar al visto en el caso de una sola capa de salida. Veamos como fluye la información por la red cuando le presentamos a la red el ejemplo μ del conjunto de entrenamiento.

Las neuronas de la capa oculta calculan sus estímulos:

$$h_j^\mu = \sum_{k=1}^N \omega_{jk} \xi_k^\mu = \bar{\omega}_j \cdot \bar{\xi}^\mu$$

Notemos que llamamos w_{jk} (minúscula) a la sinapsis entre la entrada k y la neurona j de la capa de entrada.

La salida de las neuronas de la capa oculta será entonces:

$$V_j^\mu = g_1(h_j^\mu) = g_1(\bar{\omega}_j \cdot \bar{\xi}^\mu)$$

Con los valores de las neuronas de la capa oculta podemos procesar la información en cada una de las M neuronas de la capa de salida:

$$h_i^M = \sum_{j=1}^L w_{ij} \cdot v_j^M = \bar{w}_i \cdot \bar{v}^M$$

$$\begin{aligned} \mathcal{O}_i^M &= g_2(h_i^M) = g_2(\bar{w}_i \cdot \bar{v}^M) \\ &= g_2\left(\sum_{j=1}^L w_{ij} g_1\left(\sum_{k=1}^N w_{jk} \mathcal{Z}_k^M\right)\right) \end{aligned}$$

Ahora que tengo las salidas de la red podemos comparar con el resultado correcto para el elemento μ del conjunto de entrenamiento.

$$\begin{aligned} E[\{\bar{w}\}] &= \frac{1}{2} \sum_{\mu=1}^P \sum_{i=1}^M [\mathcal{Z}_i^M - \mathcal{O}_i^M]^2 \\ &= \frac{1}{2} \sum_{\mu=1}^P \sum_{i=1}^M \left[\mathcal{Z}_i^M - g_2\left(\sum_{j=1}^L w_{ij} g_1\left(\sum_{k=1}^N w_{jk} \mathcal{Z}_k^M\right)\right) \right]^2 \end{aligned}$$

$$g_2' g_1' \frac{\partial}{\partial w} \sum_k w_{jk} \mathcal{Z}_k^M$$

Una vez que tenemos la expresión funcional de la ECM en función de todas las componentes de los vectores \overline{W}_i y \overline{w}_j , vamos a aplicar el descenso por el gradiente de una forma muy específica.

Primero corregimos los pesos sinápticos \overline{W}_i (mayúscula) entre la capa oculta y la capa de salida sin importarnos por los \overline{w}_j (minúscula).

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

donde:

$$\begin{aligned} \Delta W_{ij} &= -\eta \frac{\partial E}{\partial W_{ij}} \\ &= \eta \sum_{\mu=1}^P \cancel{\frac{\partial}{\partial}} \left[\cancel{\frac{\partial}{\partial}} \left[s_i^\mu - t_i^\mu \right] g_2'(h_i^\mu) \right] V_j^\mu \\ &= \eta \sum_{\mu=1}^P \delta_i^\mu V_j^\mu \end{aligned}$$

y

$$\delta_i^\mu = g_2'(h_i^\mu) [s_i^\mu - t_i^\mu]$$

Una vez que hemos actualizado todos los parámetros sinápticos de la última capa, pasamos a hacer lo mismo con los parámetros que van de la entrada a la capa oculta, o sea, los \underline{w}_j minúsculos. Para eso calculamos estas componentes del gradiente:

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

donde

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

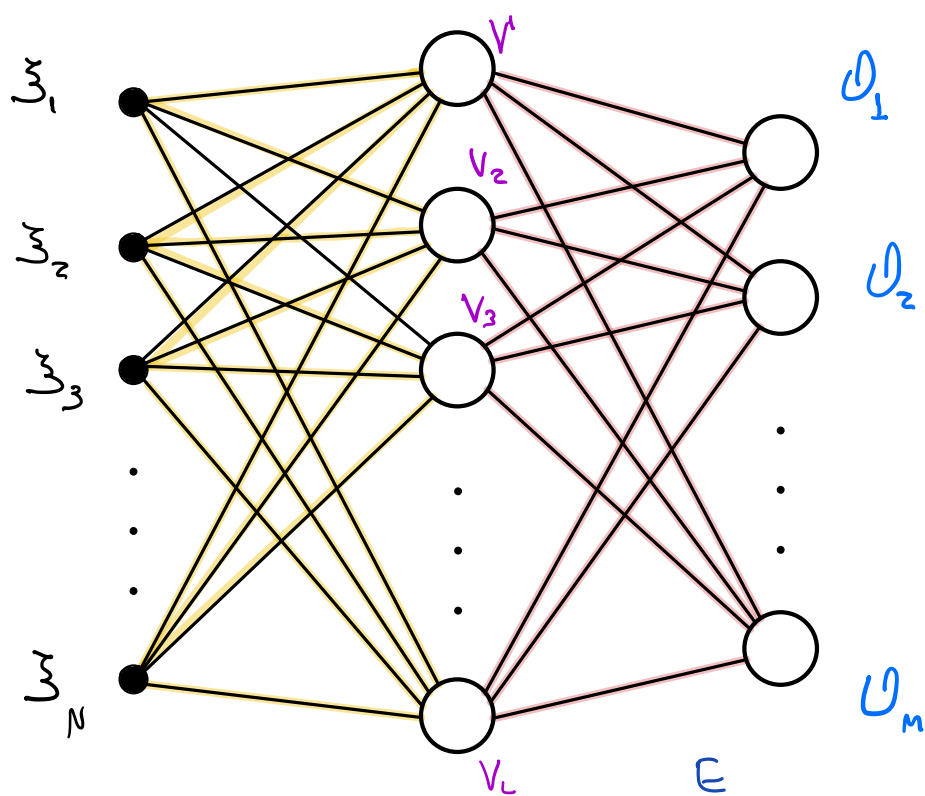
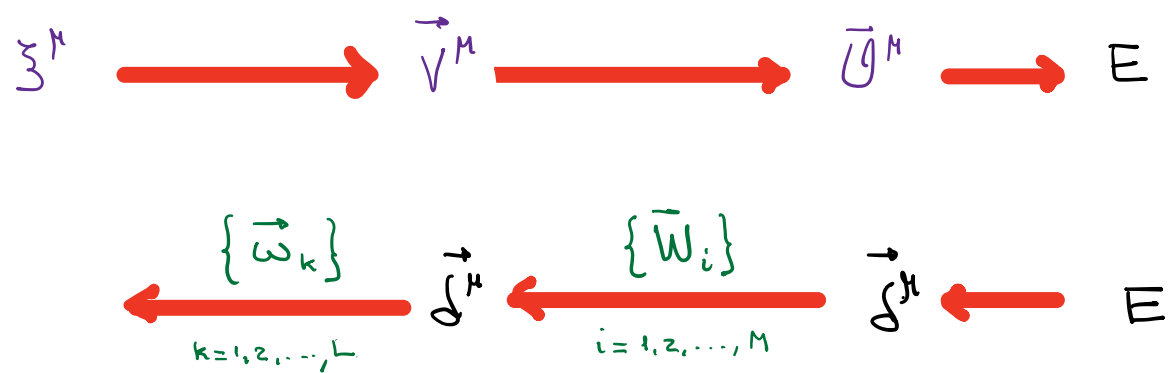
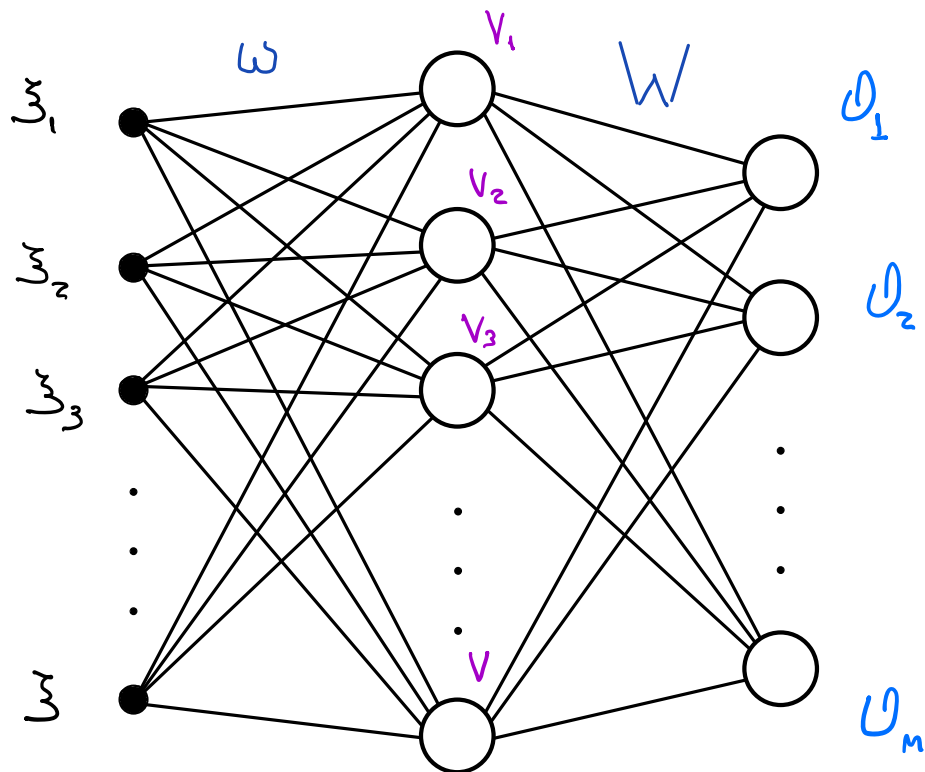
$$= \eta \sum_{\mu=1}^P [y_i^{\mu} - t_i^{\mu}] g'_2(h_i^{\mu}) w_{ij} g'_1(h_j^{\mu}) z_k^{\mu}$$

$$= \eta \sum_{\mu=1}^P \delta_i^{\mu} w_{ij} g'_1(h_j^{\mu}) z_k^{\mu}$$

$$= \eta \sum_{\mu=1}^P \delta_j^{\mu} z_k^{\mu}$$

con

$$\delta_j^{\mu} = g'_1(h_j^{\mu}) \sum_i w_{ij} \delta_i^{\mu}$$



¿Cuántos parámetros tenemos?

Este método se llama *back-propagation* o retro-propagación. Noten que le mostramos el elemento $(N \times L) + (L \times M) = L \times (N + M)$ del conjunto de entrenamiento a la red, la información viaja hacia adelante, en dirección a la salida. Con los resultados de las M salidas calculamos el error. Con el error primero calculamos las correcciones a los pesos sinápticos w entre la capa oculta y la capa de salida. Una vez actualizados estos parámetros sinápticos pasamos a corregir los pesos sinápticos w_i entre la entrada y la capa oculta.

j

El algoritmo ÉPOCA

A. Sea $\mu = 1$

B. Mientras $(\mu \leq p)$ repetimos

1. Con $\bar{\mathbf{z}}^\mu$ calculamos $\bar{\mathbf{v}}^\mu$

2. Con $\bar{\mathbf{v}}^\mu$ calculamos $\bar{\mathbf{o}}^\mu$

3. Con $\bar{\mathbf{o}}^\mu$ y $\bar{\mathbf{z}}^\mu$ calculamos E

4. Con E calculamos $\bar{\mathcal{J}}_i^\mu$

5. Con $\bar{\mathcal{J}}_i^\mu$ calculamos los $\Delta \bar{\mathbf{w}}$ y los acumulamos

$$\Delta \bar{\mathbf{w}} = \Delta \bar{\mathbf{w}} + \Delta \bar{\mathbf{w}}^\mu$$

6. Con $\bar{\mathcal{J}}_i^\mu$ calculamos los $\bar{\mathcal{J}}_j$

7. Con $\bar{\mathcal{J}}_j$ actualizamos los $\Delta \bar{\mathbf{w}}$ y los acumulamos

$$\Delta \bar{\mathbf{w}} = \Delta \bar{\mathbf{w}} + \Delta \bar{\mathbf{w}}^\mu$$

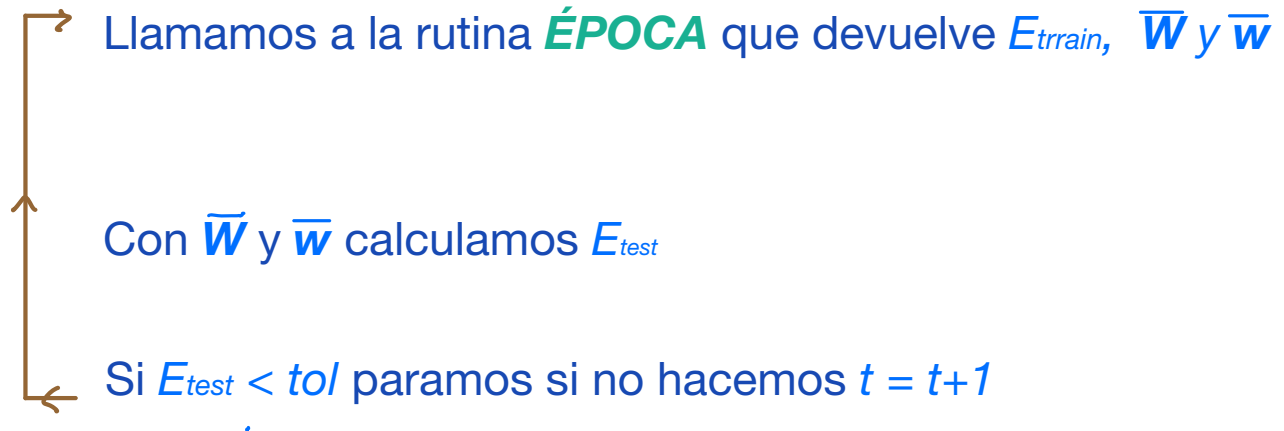
8. Pasamos al próximo ejemplo $\mu = \mu + 1$

9. Si $\mu = p$ actualizamos los $\bar{\mathbf{w}}_i$ y $\bar{\mathbf{w}}_j$ y volvemos a B

Leemos los datos:

- Conjunto de entrenamiento $(\bar{\mathbf{X}}, \bar{\mathbf{Y}})$
- Tolerancia tol
- Razón de aprendizaje η

Sea $t=1$



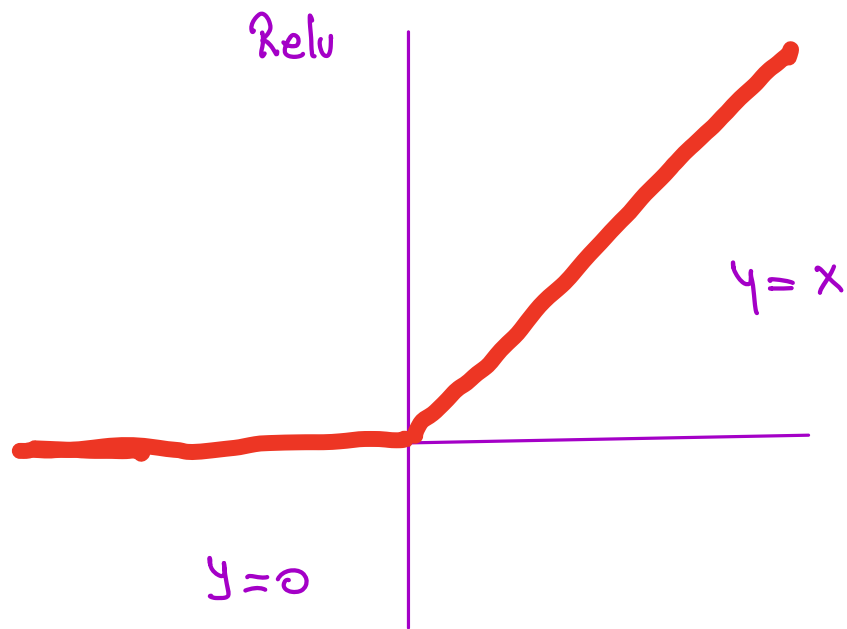
Hasta acá supusimos que actualizamos los acoplamientos sinápticos \bar{W} y \bar{w} después de mostrarle a la red todos los ejemplos del conjunto de entrenamiento, o sea después de una época. Esta forma de actualizar se conoce como *actualización en lote (batch)*.

$$\Delta W_{ij} = \Delta W_{ij}^{(1)} + \Delta W_{ij}^{(2)} + \dots + \Delta W_{ij}^{(p)}$$

$$\Delta w_{jk} = \Delta w_{jk}^{(1)} + \Delta w_{jk}^{(2)} + \dots + \Delta w_{jk}^{(p)}$$

Otra posibilidad que no hemos analizado es actualizar todos los pesos sinápticos \bar{W}_i y \bar{w}_j después de presentarle cada ejemplo del conjunto de entrenamiento a la red. Este método se denomina *actualización en línea (on line)*. Es una actualización más fina y precisa pero requiere mucho más cálculo numérico.

Veremos pronto que la forma más adecuada es una actualización en *mini-lotes (mini-batch)*, o sea, algo intermedio entre los métodos en lote y en línea.



$$\text{Relu}(x) = \begin{cases} 0 & \text{si } x < 0 \\ x & \text{si } x \geq 0 \end{cases}$$
$$= \max(0, x)$$