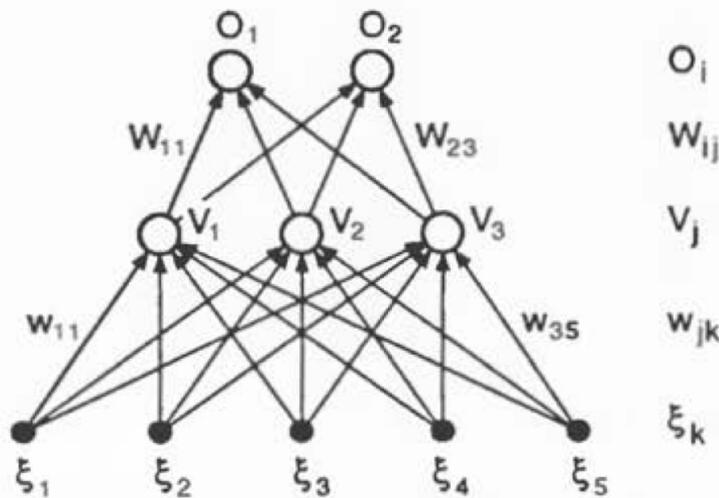


REDES MULTI CAPAS

Esta será una de las clases más importantes de este curso. Veremos como usar redes multi-capas, con $M \geq 2$.
O sea, tenemos los N elementos de entrada que definen un vector en \mathbb{R}^N

$$\vec{\xi} = (\xi_1, \xi_2, \xi_3, \dots, \xi_N),$$

Tenemos M neuronas de salida y el resto una capa oculta. Consideremos el caso de una capa oculta ($M=2$) y luego será fácil generalizarlo.
El gráfico siguiente nos da un esquema general



Las capas se organizan de abajo hacia arriba

Una vez que la red lee los N entradas ξ_k podemos calcular los estados de los neuronas V_j de la capa oculta.

$$V_j^M = g(h_j^M) = g\left(\sum_k w_{jk} \xi_k^M\right)$$

Las variables de entrada ξ_k^M pueden ser 0 o 1, +1 o -1, o puede tomar valores continuos, dependiendo de la forma de $g(h)$.

Ahora con los V_j^M podemos calcular las neuronas de salida

$$\begin{aligned} O_i^M &= g(h_i^M) \\ &= g\left(\sum_j w_{ij} V_j^M\right) \\ &= g\left(\sum_j w_{ij} g\left(\sum_k w_{jk} \xi_k^M\right)\right) \end{aligned}$$

Noteu que ahora no podemos referir la red en subredes simples con 1 neurona de salida.

La función costo es la misma

$$E(\vec{w}, \vec{W}) = \frac{1}{2} \sum_{j\mu} (y_i^\mu - O_i^\mu)^2$$
$$= \frac{1}{2} \sum_{j\mu} \left[y_i^\mu - g \left(\sum_j W_{ij} g \left(\sum_k W_{jk} z_k^\mu \right) \right) \right]^2$$

Si tenemos N entradas

L neuronas ocultas en la 1ª capa

M neuronas de salida en la 2ª capa

tendremos

$N \times L$ elementos w_{jk}

$L \times M$ elementos W_{ij}

Para generalizar el método de descenso por el gradiente para encontrar un conjunto (\vec{w}, \vec{W}) tal que $E(\vec{w}, \vec{W}) = 0$ debemos encontrar el gradiente de E , que será un campo vectorial

$$\nabla E : \mathbb{R}^{L(N+M)} \rightarrow \mathbb{R}^{L(N+M)}$$

Como ven, es un problema al menos complicado

$$\nabla E = \left(\underbrace{\frac{\partial E}{\partial w_{11}}, \frac{\partial E}{\partial w_{12}}, \dots, \frac{\partial E}{\partial w_{N1}}, \frac{\partial E}{\partial w_{11}}, \dots, \frac{\partial E}{\partial w_{LM}}}_{N \times L} \right)$$

$$N \times L + L \times M = L(N+M) \text{ sinapsis!}$$

Vamos a introducir ahora un algoritmo muy famoso e importante que acompaña a todas las redes feed-forward, llamado **BACK-PROPAGATION**

BACK-PROPAGATION

El algoritmo nos brinda un método para cambiar los $L(N+M)$ pesos sinápticos en cualquier red feed-forward multi capas. La base es el método de descenso por el gradiente que vimos para el perceptron simple con salida lineal o no-lineal (continuos).

| | |
|-------------------------------|------|
| Bryson y Ho | 1969 |
| Werbos | 1974 |
| Parker | 1985 |
| Hinton, Rumelhart y Williams, | 1985 |
| Le Cun | 1985 |

mejor enfoque

Tenemos P relaciones input-output etiquetadas

$$\bar{X}^k \rightarrow \bar{Y}^k$$

donde $\bar{O}^k = \bar{Y}^k$ es la redida esperada cuando el input es \bar{X}^k .

Etiquetamos todas las sinapsis $\{W_{jk}\}$ y $\{W_{ij}\}$ de forma estocástica (como queremos).

Le presentamos a la red los P patrones del conjunto de entrenamientos \bar{X}^k uno tras otro, en forma secuencial o estocástica y acumulamos el costo. Luego hacemos descenso por el gradiente para corregir los W_{ij}

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

$$\text{con } \Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}}$$

$$= +\eta \sum_k [Y_i^k - O_i^k] g'(h_i^k) V_j^k$$

$$\Delta W_{ij} = \eta \sum_{\mu} \delta_i^{\mu} v_j^{\mu}$$

$$\text{con } \delta_i^{\mu} = g'(h_i^{\mu}) (z_i^{\mu} - 0_i^{\mu})$$

Una vez que hemos actualizado los W_{ij} pasamos a la primera capa de neuronas.

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

con

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

$$= -\eta \sum_{\mu} \frac{\partial E}{\partial v_j^{\mu}} \cdot \frac{\partial v_j^{\mu}}{\partial w_{jk}}$$

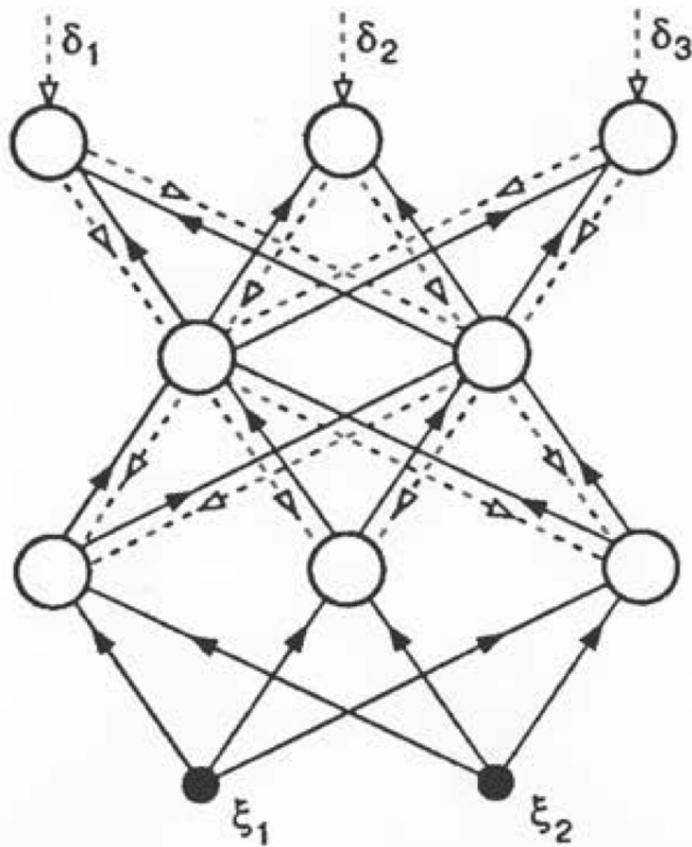
$$= \eta \sum_{\mu i} [z_i^{\mu} - 0_i^{\mu}] g'(h_i^{\mu}) w_{ij} g'(h_j^{\mu}) z_k^{\mu}$$

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

$$\Delta w_{jk} = \eta \sum_{\mu} \delta_j^{\mu} z_k^{\mu}$$

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

Se llama back-propagation porque la información viaja hacia adelante y las correcciones hacia atrás.



Noten que hay una forma general

$$\Delta w_{pq} = \eta \sum_{\mu} f_{salida} \times V_{input}$$

entre dos capas subsiguientes p y q .

Podemos actualizar

* después de cada patrón μ : on line

* después de p patrones : batch.

Podemos mezclar diferentes tipos de funciones de activación.

Podemos tener muchos copos.

El conjunto de entrenamiento puede tener diferentes tamaños

Δ o luego de la actualización, cada vez que terminamos de actualizar Δ vez cada uno de los p elementos de conjunto de entrenamiento (una época) monitoreamos la función la evolución de la **función costo del entrenamiento**

Si nos hemos guardado un conjunto de q elementos etiquetados correctamente podemos ver como evoluciona **la función costo de testeo**

Resumen: Consideremos que tenemos muchos capas $m = 0, 1, 2, \dots, M$. 0 corresponde a la entrada que no es una capa.

1. Elegimos todos los pesos sinápticos entre todos los capas y entre la entrada y la 1ª capa. los escogemos aleatorios y pequeños.

2. Tenemos un patrón \sum_k^M y lo ponemos en la entrada

$$V_k^{(0)} = \sum_k^M \quad \forall k$$

3. Propagamos la señal hacia adelante

$$V_i^{(m)} = g(h_i^{(m)}) = g\left(\sum_j \omega_{ij}^{(m)} V_j^{(m-1)}\right)$$

para cada i de la capa i -ésima y
hacia atrás a la última capa $V_i^{(M)} = O_i$

4. Calculamos los deltas de la salida

$$\delta_i^{(m)} = g' \left(h_i^{(m)} \right) \left[y_i^{(m)} - v_i^{(m)} \right]$$

comparando lo obtenido con lo deseado.

5. Calculamos los δ de las capas anteriores, en orden decreciente

$$\delta_i^{(m-1)} = g' \left(h_i^{(m-1)} \right) \sum_j w_{ji}^{(m)} \delta_j^{(m)}$$

para $m = \bar{m}-1, \bar{m}-2, \dots, 2$ para cada neurona de cada capa.

6. Aplicamos la regla

$$w_{ij}^{(m) \text{ nuevo}} = w_{ij}^{(m) \text{ viejo}} + \Delta w_{ij}^{(m)}$$

$$\text{con } \Delta w_{ij}^{(m)} = \eta \delta_i^{(m)} v_j^{(m-1)}$$

7. Volvemos al punto (2)

MOMENTO

$$\Delta W_{ij}^{(m)}(t+1) = -\eta \frac{\partial E}{\partial W_{ij}^{(m)}} + \alpha \Delta W_{ij}^{(m)}(t)$$

Quando o gradiente for uma soma muito pequena

$$\Delta W_{ij}^{(m)}(t+1) \approx \Delta W_{ij}^{(m)}(t)$$

$$\Delta W_{ij}^{(m)}(t+1) (1-\alpha) \approx -\eta \frac{\partial E}{\partial W_{ij}^{(m)}}$$

$$\Delta W_{ij}^{(m)} \approx \underbrace{-\frac{\eta}{(1-\alpha)}}_{\text{efectivo}} \frac{\partial E}{\partial W_{ij}^{(m)}}$$

PARAMETRO ADAPTATIVO

$$\Delta\eta = \begin{cases} +\alpha & \text{si } \Delta E < 0 \text{ consistentemente} \\ -\beta\eta & \text{si } \Delta E > 0 \\ 0 & \text{si no} \end{cases}$$

k pasos

OTROS MÉTODOS DE APROXIMACIÓN

no lo voy hacer este en la
rección 6.2 del capítulo 6 del
libro de Hertz et al.

APLICACIONES

XOR

PARITY FUNCTION

⋮

TEOREMA UNIVERSAL DE LA APROXIMACIÓN

Una red feed-forward con una capa oculta ($M=2$) que contiene un número finito de neuronas puede aproximar funciones continuas sobre un subconjunto compacto de \mathbb{R}^n , para funciones de activación con buenas propiedades.

1989 George Cybenko (sigmoideas)

En 2017 Lu et al probaron esto para redes profundas. $(4+n)$ con ReLU. EXPLICAR