

# Análisis Matemático para Inteligencia Artificial

Verónica Pastor (vpastor@fi.uba.ar),  
Martín Errázquin (merrazquin@fi.uba.ar)

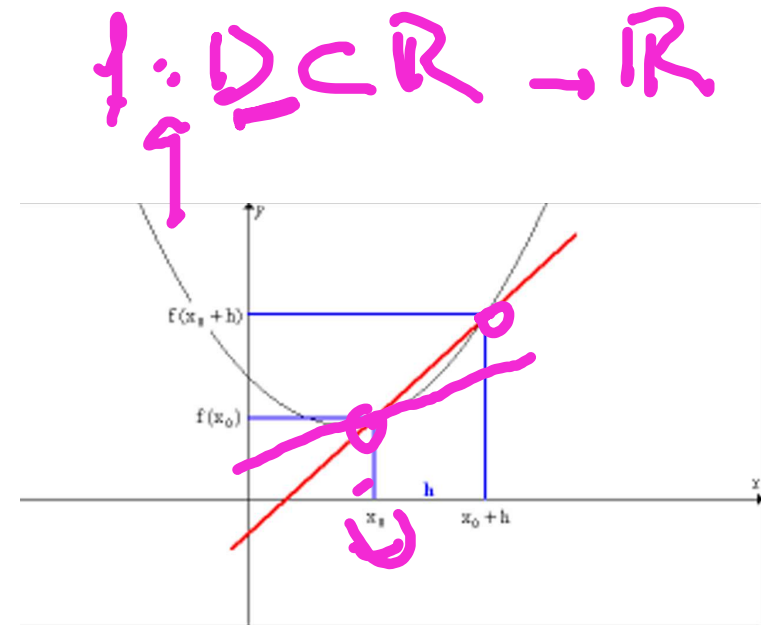
Especialización en Inteligencia Artificial

31/03/2023

# Análisis Matemático - Optimización

# Repaso

- 1 En los videos de repaso definimos funciones de cuyo dominio y codominio eran los reales, la gráfica de la función se representa en  $\mathbb{R}^2$ .
- 2 Toda función  $f$  describe el cambio de una magnitud (v. dependiente) en términos de otra (v. independiente), cuando esta variable se mueve en cierto intervalo  $[x_0, x_0 + h]$  la variación total se mide como  $f(x_0 + h) - f(x_0)$ .
- 3 Mientras que la variación media es  $\frac{f(x_0 + h) - f(x_0)}{(x_0 + h) - x_0}$ . Geométricamente, podemos ver la variación media como la pendiente de la recta secante.
- 4 Cuando hacemos que  $h \rightarrow 0$ , ...



... esto nos conduce a la definición de derivada de  $f$  en

$x_0$ :

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

# Clasificación de funciones

Dada  $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ .

- Si  $m = 1$  diremos que es una función

- **escalar**, si  $n = 1$ ,

- **campo escalar**,  $n > 1$ .

$$f: D \subset \mathbb{R} \rightarrow \mathbb{R} \quad \checkmark$$

$$f: D \subset \mathbb{R}^2 \rightarrow \mathbb{R} \quad (*) \text{ ejemplo}$$

- Si  $m > 1$  diremos que es una función

- **vectorial**, si  $n = 1$ ,

- **campo vectorial**,  $n > 1$ .

$$f: D \subset \mathbb{R} \rightarrow \mathbb{R}^2 \quad f(t) = \begin{pmatrix} t^2 \\ e^t \end{pmatrix} \text{ derivado}$$

$$f_1(t) \quad f_2(t)$$

$$f'(t) = (2t, e^t)$$

**Conjuntos de Nivel** Dada  $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  el conjunto de nivel  $k$  de  $f$ ,  $L_k \subset \mathbb{R}^n$ , definido por:

$$L_k = \{x \in \mathbb{R}^n / x \in D \wedge f(x) = k\}$$

La representación geométrica de  $L_k$  se obtiene identificando gráficamente los puntos del dominio de la función para los cuales el valor de  $f$  es igual a  $k$ , para graficar no es necesario agregar un eje.

$$f(x, y) = x^2 + y^2$$

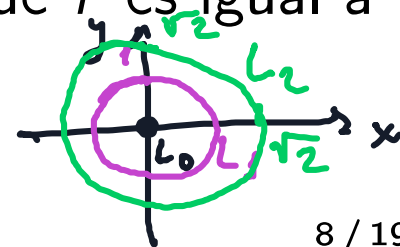
$$\text{Dom } f: \mathbb{R}^2$$

$$L_0 = \{x \in \mathbb{R}^2 : x \in \mathbb{R}^2 \wedge x^2 + y^2 = 0\}$$

$$L_1 = \{x \in \mathbb{R}^2 : x^2 + y^2 = 1\}$$

$$L_2 = \{x \in \mathbb{R}^2 : x^2 + y^2 = 2\} \rightarrow r^2 \rightarrow r = \sqrt{2}$$

$k < 0?$   
 $n = 2$  tiene



# Derivando campos ...

- escalares: Sea  $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $(x_1, \dots, x_n)^T \mapsto f((x_1, \dots, x_n)^T)$ , se definen las **derivadas parciales** como:

$$\frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h}$$

Handwritten notes:  $x_0 - h, x_0, x_0 + h$  with arrows;  $\vec{n} = (0, 1)$  with arrows pointing to the  $x_2$  term in the denominator.

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h}$$

Se define el **gradiente** como:  $\nabla f = \left( \frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n} \right)$ .

Importante: El gradiente apunta en la dirección de máximo crecimiento.

- vectoriales: Sea  $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $(x_1, \dots, x_n)^T \mapsto (f_1((x_1, \dots, x_n)^T), \dots, f_m((x_1, \dots, x_n)^T))$ , se define el **jacobiano** como:

Handwritten notes:

$$f_i((x_1, \dots, x_n)) = r_{i0} \in \mathbb{R}$$

$$f_i: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$$

$$i = 1, \dots, m$$

$$J_f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Handwritten notes:  $\rightarrow \nabla f_1$  and  $\rightarrow \nabla f_m$  with arrows pointing to the first and last rows of the Jacobian matrix.

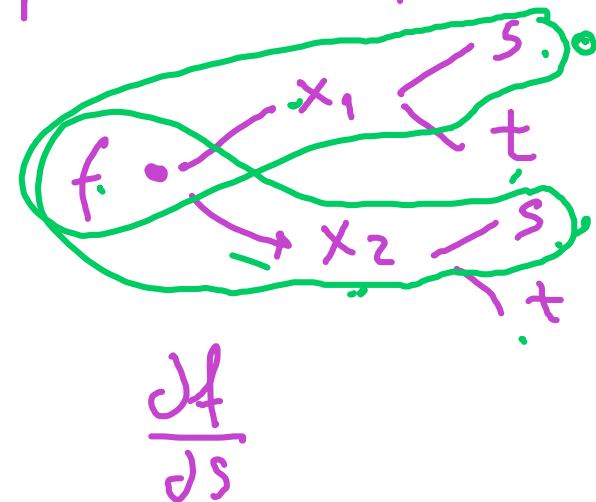
# Regla de la Cadena en forma matricial

Sea  $f(x_1(s, t), x_2(s, t))$

Para derivar una composición de funciones

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s}$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$



Y luego

$$\frac{df}{d(s, t)} = \frac{df}{dx} \frac{dx}{d(s, t)} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix} \begin{matrix} \sim \nabla x_1 \\ \sim \nabla x_2 \end{matrix}$$

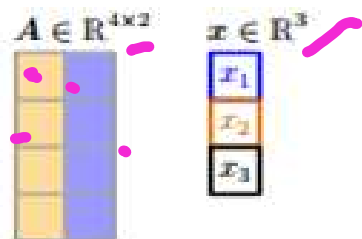
Recordemos reglas de derivación:

- $\frac{\partial (f+g)(s)}{\partial s} = \frac{\partial f}{\partial s} + \frac{\partial g}{\partial s}$
- $\frac{\partial (fg)(s)}{\partial s} = \frac{\partial f}{\partial s} g(s) + f(s) \frac{\partial g}{\partial s}$

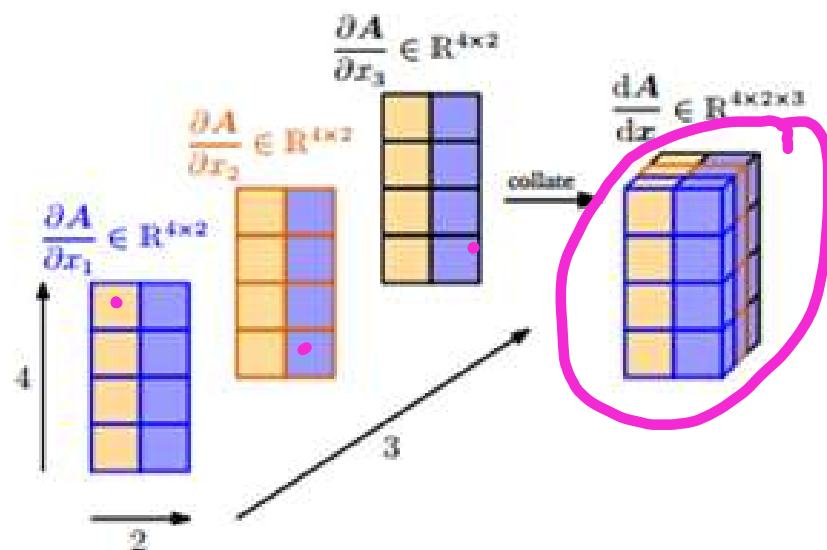
$$\frac{d}{ds} \left( \frac{f}{g} \right) (s) = \frac{\frac{df}{ds} g(s) - f(s) \frac{dg}{ds}}{(g(s))^2}$$

•  $\frac{d}{ds} (f \cdot s^{-1})$

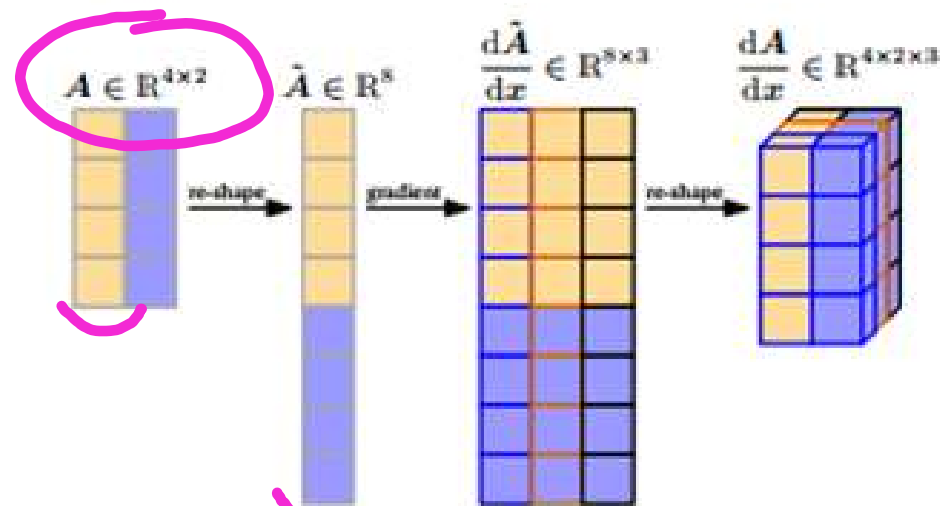
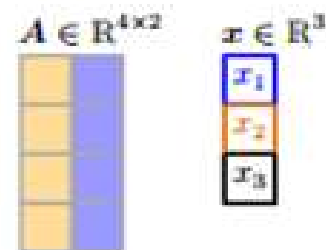
# Derivada de matrices



Partial derivatives:



(a) Approach 1: We compute the partial derivative  $\frac{\partial A}{\partial x_1}$ ,  $\frac{\partial A}{\partial x_2}$ ,  $\frac{\partial A}{\partial x_3}$ , each of which is a  $4 \times 2$  matrix, and collate them in a  $4 \times 2 \times 3$  tensor.



(b) Approach 2: We re-shape (flatten)  $A \in \mathbb{R}^{4 \times 2}$  into a vector  $\tilde{A} \in \mathbb{R}^8$ . Then, we compute the gradient  $\frac{d\tilde{A}}{dx} \in \mathbb{R}^{8 \times 3}$ . We obtain the gradient tensor by re-shaping this gradient as illustrated above.

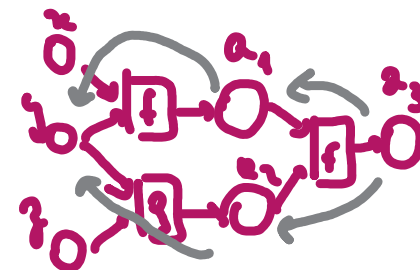
# Diferenciación Automática

Sean, para una función  $f$ :

- $x_1, \dots, x_d$  las variables de entrada
- $x_{d+1}, \dots, x_{D-1}$  las variables intermedias
- $x_D$  la variable de salida
- $g_i$  funciones elementales
- $Hij(x_i)$  el conjunto de nodos hijos de cada  $x_i$

Así queda definido un **grafo de cómputo**. Recordando que  $f = D$ , tenemos que  $\frac{\partial f}{\partial x_D} = 1$ . Para las otras variables  $x_i$  aplicamos la regla de la cadena:

$$\frac{\partial f}{\partial x_i} = \sum_{x_j \in H_{ij}(x_i)} \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial x_i} = \sum_{x_j \in H_{ij}(x_i)} \frac{\partial f}{\partial g_j} \frac{\partial x_j}{\partial x_i}$$



- La diferenciación automática se puede utilizar siempre que la función pueda representarse como un grafo de cómputo.
- La gran ganancia de este mecanismo está en que cada función sólo precisa saber cómo derivarse a sí misma, permitiendo OOP.

```
class prod:
    def f(self, x, y):
        return x * y
    def df(self, x, y):
        return (y, x)
```

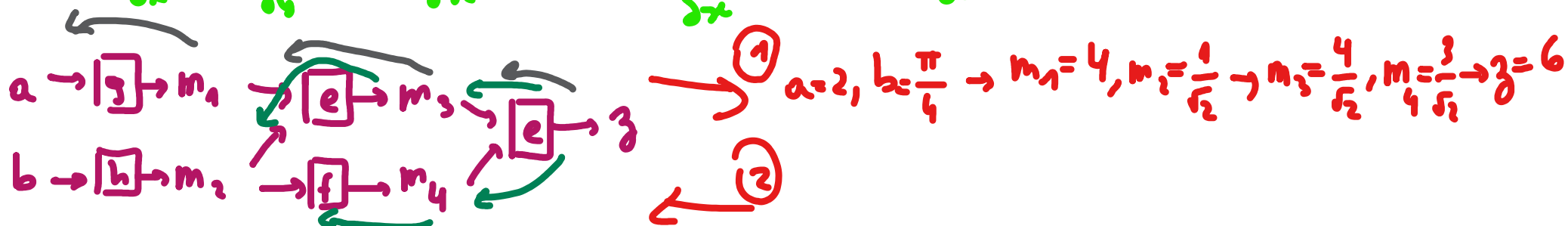


# Diferenciación Automática: ejemplo

Sean  $e(x, y) = xy$ ,  $f(x) = 3x$ ,  $g(x) = x^2$ ,  $h(x) = \sin(x)$

↳ quiero  $a=2, b=\frac{\pi}{4}$

$$\frac{\partial e}{\partial x} = y, \frac{\partial e}{\partial y} = x \quad \frac{\partial f}{\partial x} = 3 \quad \frac{\partial g}{\partial x} = 2x \quad \frac{\partial h}{\partial x} = \cos(x)$$



$$\left[ \frac{\partial z}{\partial a} = \frac{\partial z}{\partial m_3} \cdot \frac{\partial m_3}{\partial m_1} \cdot \frac{\partial m_1}{\partial a} = m_4 \cdot m_2 \cdot 2a = \frac{3}{\sqrt{2}} \cdot \frac{1}{\sqrt{2}} \cdot 4 = \frac{12}{2} = 6 \right]$$

$$\left[ \frac{\partial z}{\partial b} = \frac{\partial z}{\partial m_2} \cdot \frac{\partial m_2}{\partial b} = \left( \frac{\partial z}{\partial m_3} \cdot \frac{\partial m_3}{\partial m_2} + \frac{\partial z}{\partial m_4} \cdot \frac{\partial m_4}{\partial m_2} \right) \cdot \frac{\partial m_2}{\partial b} = (m_4 \cdot m_1 + m_3 \cdot 3) \cdot \cos(b) \right]$$

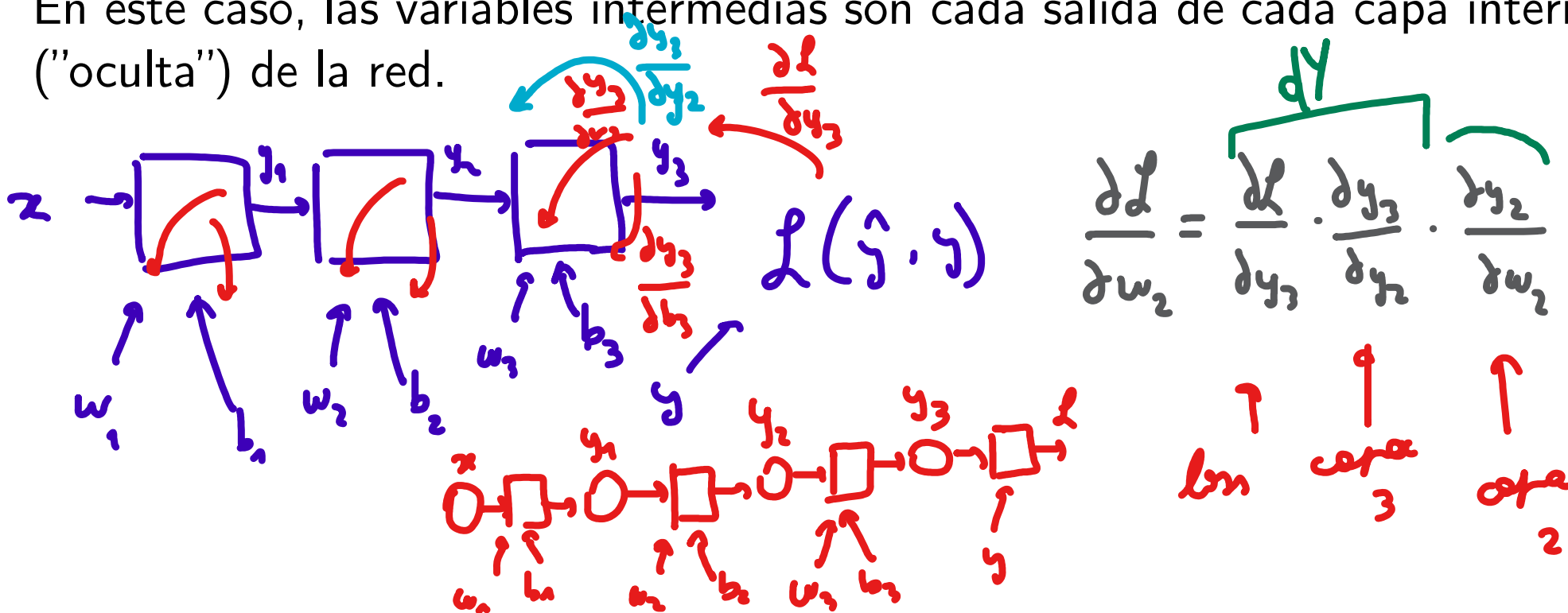
$$= \left( \frac{3}{\sqrt{2}} \cdot 4 + \frac{4}{\sqrt{2}} \cdot 3 \right) \cdot \frac{1}{\sqrt{2}} = \frac{24}{\sqrt{2}} \cdot \frac{1}{\sqrt{2}} = 12$$

# Backpropagation

¿Dónde se aplica la diferenciación automática? En **Backpropagation** (o simplemente Backprop), el algoritmo utilizado para entrenar redes neuronales.

¿Qué función cumple? La de computar las derivadas de la función de error/costo respecto de *cada* parámetro de la red neuronal.

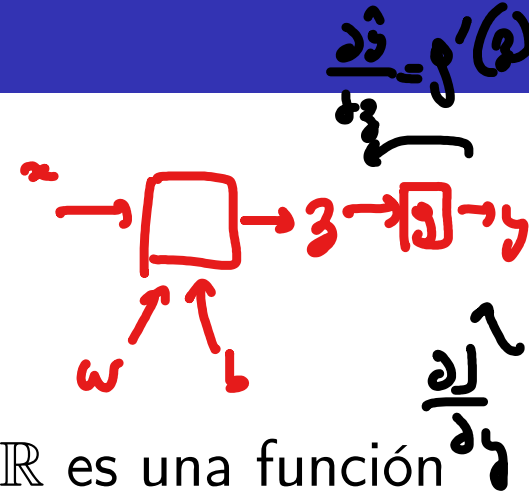
En este caso, las variables intermedias son cada salida de cada capa interna ("oculta") de la red.



# Redes neuronales (I)

Un perceptrón/neurona es un estimador de la forma:

$$\hat{y} = g(w \cdot x + b)$$



donde en su forma más simple  $x, y, w, b \in \mathbb{R}$  y  $g : \mathbb{R} \rightarrow \mathbb{R}$  es una función no lineal como puede ser la sigmoidea  $\sigma(z) = \frac{1}{1+e^{-z}}$ .

Si se define la función  $J(W, b)$  de error respecto de los parámetros  $W$  y  $b$  se puede comprobar que, definiendo  $z = \underline{w \cdot x + b}$  y suponiendo conocido

$$\underbrace{\frac{dJ}{d\hat{y}}}_{dY} \in \mathbb{R}:$$

$$\frac{\partial \hat{y}}{\partial z} = g'(z)$$

$$\frac{\partial J}{\partial W} = \frac{dJ}{d\hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial W} = \overset{\mathbb{R}}{dY} \cdot \overset{\mathbb{R}}{g'(z)} \cdot \overset{\mathbb{R}}{\hat{x}} \in \mathbb{R} \quad \checkmark$$

$$\frac{\partial J}{\partial b} = \frac{dJ}{d\hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b} = dY \cdot g'(z) \cdot 1 \in \mathbb{R} \quad \checkmark$$

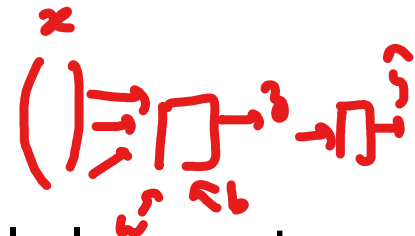
# Redes neuronales (II)

Si ahora consideramos múltiples entradas, es decir  $x \in \mathbb{R}^n$ ,  $W \in \mathbb{R}^{1 \times n}$ :

$$x = (x_1, \dots, x_n)^T$$

$$W = (w_1, \dots, w_n)$$

$$\hat{y} = g(\underbrace{W}_{1 \times n} \cdot \underbrace{x}_n + b)$$



Entonces ahora para cada elemento de  $W = (w_1, \dots, w_n)$  vale lo anterior, y por tanto se puede comprobar que

$$\frac{\partial J}{\partial W} = \nabla_J(W) = (\underbrace{\overbrace{dY \cdot g'(z)}^{\mathbb{R}} \cdot x_1}_{dw_1}, \dots, \underbrace{\overbrace{dY \cdot g'(z)}^{\mathbb{R}} \cdot x_n}_{dw_n}) = \overbrace{dY}^1 \cdot \overbrace{g'(z)}^1 \cdot \underbrace{x^T}_{1 \times n} \in \mathbb{R}^{1 \times n} \checkmark$$

$$\frac{\partial J}{\partial b} = dY \cdot g'(z) \cdot 1 \quad \checkmark$$

# Redes neuronales (III)

Una capa en una red neuronal se define como un vector de  $k$  neuronas en paralelo. Una propiedad atractiva de este formato es que se puede considerar a la salida de una capa  $y \in \mathbb{R}^k$  como simplemente el  $x$  de la capa siguiente. Por convención (y eficiencia computacional) se suele utilizar la misma no-linealidad  $g$  para todas las neuronas de la capa.

Nuevamente tenemos:

$$\hat{y} = g(\underbrace{W}_{k \times n} \cdot \underbrace{\hat{x}}_{n \times 1} + \underbrace{b}_{k \times 1})$$

donde  $x \in \mathbb{R}^n$ ,  $W \in \mathbb{R}^{k \times n}$ ,  $b \in \mathbb{R}^k$  y se conviene  $g(z) = \begin{pmatrix} g(z_1) \\ \vdots \\ g(z_k) \end{pmatrix}$

¿Y ahora cómo se calculan las derivadas para  $W$  y  $b$ ?

# Redes neuronales (IV)

En el caso de  $b$  es simple:

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b} = \begin{pmatrix} dY_1 \\ \vdots \\ dY_k \end{pmatrix} \odot \begin{pmatrix} g'(z_1) \\ \vdots \\ g'(z_k) \end{pmatrix} \odot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \underbrace{dY}_{k \times 1} \odot \underbrace{g'(z)}_{k \times 1} \in \mathbb{R}^k \checkmark$$

prod. elemento a elemento  
↓

Ahora para cada elemento de  $W$  tenemos:

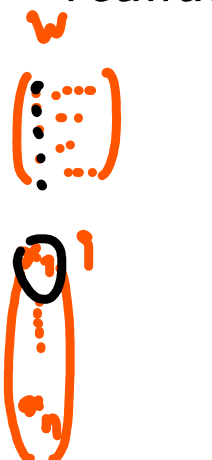
$$\frac{\partial J}{\partial W} = \begin{pmatrix} \frac{\partial J}{\partial W_{1,1}} & \cdots & \frac{\partial J}{\partial W_{1,n}} \\ \frac{\partial J}{\partial W_{2,1}} & \cdots & \frac{\partial J}{\partial W_{2,n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial J}{\partial W_{k,1}} & \cdots & \frac{\partial J}{\partial W_{k,n}} \end{pmatrix} = \begin{pmatrix} \nabla_J(W_{1,:}) \\ \vdots \\ \nabla_J(W_{k,:}) \end{pmatrix} = \begin{pmatrix} dY_1 \cdot g'(z_1) \cdot x^T \\ \vdots \\ dY_k \cdot g'(z_k) \cdot x^T \end{pmatrix} =$$

(dY \* dg) @ X.T

$$= \begin{pmatrix} dY_1 \\ \vdots \\ dY_k \end{pmatrix} \odot \begin{pmatrix} g'(z_1) \\ \vdots \\ g'(z_k) \end{pmatrix} \cdot x^T = \underbrace{dY \odot g'(z)}_{k \times 1} \cdot \underbrace{x^T}_{1 \times n} \in \mathbb{R}^{k \times n} \checkmark$$

# Redes neuronales (V): Backpropagation

¿Cómo se encadena esto? Nosotros estamos dando por conocida la derivada del error respecto de la salida de la capa,  $dY = \frac{dJ}{d\hat{y}}$ , pero en realidad no tenemos idea si estamos en una capa intermedia o no.



$$\frac{\partial \hat{y}}{\partial x_i} = \sum_{j=1}^k \underbrace{dY_j}_{\frac{\partial J}{\partial \hat{y}_j}} \cdot \underbrace{g'(z_j)}_{\frac{\partial z_j}{\partial x_i}} \cdot \underbrace{W_{j,i}}_{\frac{\partial z_j}{\partial x_i}} = \left\langle \begin{pmatrix} dY_1 \cdot g'(z_1) \\ \vdots \\ dY_k \cdot g'(z_k) \end{pmatrix}, \begin{pmatrix} W_{1,i} \\ \vdots \\ W_{k,i} \end{pmatrix} \right\rangle =$$

$\langle x, y \rangle = x^T y = y^T x$

$$= \underbrace{\langle dY \odot g'(z), W_{:,i} \rangle}_{k \times 1} = \underbrace{W_{i,:}^T}_{1 \times k} \cdot \underbrace{(dY \odot g'(z))}_{k \times 1} \in \mathbb{R}$$

En forma vectorizada:

$$dX = \frac{\partial J}{\partial x} = \begin{pmatrix} \frac{\partial J}{\partial x_1} \\ \vdots \\ \frac{\partial J}{\partial x_n} \end{pmatrix} = \begin{pmatrix} W_{1,:}^T \cdot dY \odot g'(z) \\ \vdots \\ W_{n,:}^T \cdot dY \odot g'(z) \end{pmatrix} = \underbrace{W^T}_{n \times k} \cdot \underbrace{(dY \odot g'(z))}_{k \times 1} \in \mathbb{R}^n$$

Y ese  $dX$  no es otra cosa que el  $dY$  de la capa anterior!