

# Análisis Matemático para Inteligencia Artificial

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

Especialización en Inteligencia Artificial

1/4/2022

# Repaso

En las clases de repaso definimos funciones de cuyo dominio y codominio eran los reales, la gráfica de la función se representa en  $\mathbb{R}^2$ .

Además, dijimos que 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  $f(x_0 + h) - f(x_0)$ , 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, pero cuando hacemos que  $h \rightarrow 0$ ,



Esto nos conduce a la definición de derivada de  $f$  en  $x_0$ :

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

Handwritten orange annotations: A circle around the numerator  $f(x_0 + h) - f(x_0)$ , and a bracket under the denominator  $h$  with the text  $x_0 + h - x_0$  written below it.

# Clasificación de funciones

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

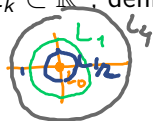
- Si  $m = 1$  diremos que es una función

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

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

- **vectorial**, si  $n = 1$ ,  $f(t) = (t^2, 2t)$  *t. parámetro*
- **campo vectorial**,  $n > 1$ .  $f(x, y) = (x^2, 0, xy)$   $f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

**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:



$$L_k = \{x \in \mathbb{R}^n / x \in D \wedge f(x) = k\} \quad \text{Dom}(f) = \mathbb{R}^2$$

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.

$$L_{-1}: x^2 + y^2 = -1 \quad (k \geq 0) \quad L_0: x^2 + y^2 = 0 \quad L_1: x^2 + y^2 = 1 \quad L_4: x^2 + y^2 = 4$$

# 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:

$$\frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$
$$\frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(\overbrace{x_1 + h}^{x_0 + h}, \overbrace{x_2, \dots, x_n}^{x_0}) - f(x_1, x_2, \dots, x_n)}{h}$$
$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h}$$

- 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 **gradiente** como:

$$\nabla f = \left( \frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right)$$
$$\left[ \begin{array}{c} \frac{\partial f_1}{\partial x_1} \quad \dots \quad \frac{\partial f_1}{\partial x_n} \\ \vdots \\ \frac{\partial f_m}{\partial x_1} \quad \dots \quad \frac{\partial f_m}{\partial x_n} \end{array} \right] = J$$

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

$|\nabla f|$  mide la pendiente máx

# Regla de la Cadena en forma matricial

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



$$x_1 = f_1(s, t) \\ x_2 = f_2(s, t)$$

$$f(s, t) = (st^2, st)$$

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

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

Y luego

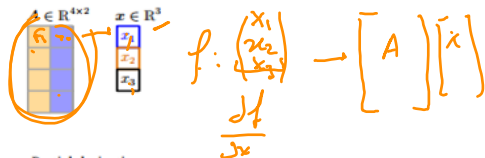
$$\frac{df}{d(s, t)} = \frac{df}{dx} \frac{dx}{d(s, t)} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial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

Recordemos reglas de derivación:

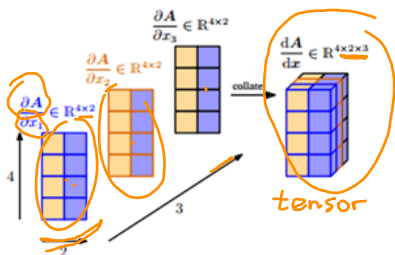
- $\frac{\partial (f+g)(s)}{\partial s} = \frac{\partial f}{\partial s} + \frac{\partial g}{\partial s}$
- $\frac{\partial (fg)(s)}{\partial s} = \frac{\partial f}{\partial s} \underline{g(s)} + \underline{f(s)} \frac{\partial g}{\partial s}$

$$(f \cdot g)(s) =$$

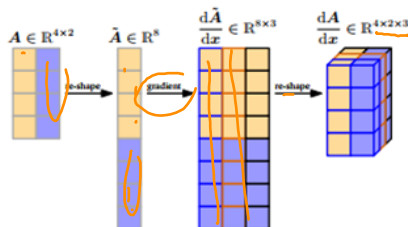
# Derivada de matrices



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.

# Matriz Hessiana

La **matriz Hessiana** es aquella cuyas derivadas de orden 2 de  $f$  respecto a  $x \in \mathbb{R}^n$  se ubican:

$\left(\frac{\partial f}{\partial x}\right)$  *deriva 2 veces f*

*deriva 2 veces respecto a  $x_1$*

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Sea  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$   $f(x,y) = 2xy^2$

$$\frac{\partial f}{\partial x} = 2y^2 \quad \begin{cases} \frac{\partial^2 f}{\partial x^2} = 0 \\ \frac{\partial^2 f}{\partial y \partial x} = 4y \end{cases}$$

$$\frac{\partial f}{\partial y} = 4xy \quad \begin{cases} \frac{\partial^2 f}{\partial x \partial y} = 4y \\ \frac{\partial^2 f}{\partial y^2} = 4x \end{cases}$$

$\nabla f$

$$S^T x = \tilde{x}$$

$$\underline{x^T H x} = x^T S D S^T x = (S^T x)^T D (S^T x) = \underline{\tilde{x}^T D \tilde{x}} \quad \begin{matrix} H = S D S^T \\ \rightarrow \text{F. cuadrática} \end{matrix}$$

$$H = \begin{bmatrix} 0 & 4y \\ 4y & 4x \end{bmatrix}$$

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x} \quad H^T = \begin{bmatrix} 0 & 4y \\ 4y & 4x \end{bmatrix}$$

$H = H^T$ , es simétrica

$$\exists S: S^T H S = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$

Cde  
variable

# Aplicación: Polinomio de Taylor

Sea  $f$  un campo escalar  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , asumiendo que posee derivadas parciales de todo orden en un entorno de un punto  $a \in \mathbb{R}^n$ , se define el **polinomio de Taylor** de grado  $k$ :

$$\begin{aligned} P_k(x) = & \underbrace{f(a)} + \sum_{i=1}^n \underbrace{\frac{\partial f}{\partial x_i}(a)} \underbrace{(x_i - a_i)} + \frac{1}{2!} \sum_{i,j=1}^n \underbrace{\frac{\partial^2 f}{\partial x_i \partial x_j}(a)} \underbrace{(x_i - a_i)(x_j - a_j)} + \dots + \\ & \dots + \frac{1}{k!} \sum_{\text{indices}} \frac{\partial^k f}{\partial x_{i_1} \dots \partial x_{i_k}}(a) \underbrace{(x_{i_1} - a_{i_1}) \dots (x_{i_k} - a_{i_k})} + \mathcal{O}(k+1) \end{aligned}$$

↳ todas las comb. posibles de  $k$  elementos en el conjunto  $x_1, \dots, x_n$



$T(a+\varepsilon)$  fra  $\varepsilon$

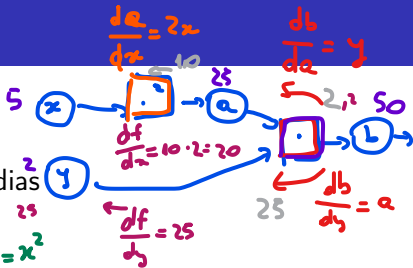




# Diferenciación Automática

Sean, para una función  $f$ :  $b = x^2 \cdot y$

- $x_1, \dots, x_d$  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:

$$\frac{\partial f}{\partial x_i} = \sum_{x_j \in Hij(x_i)} \left[ \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial x_i} \right] = \sum_{x_j \in Hij(x_i)} \frac{\partial f}{\partial g_j} \frac{\partial x_j}{\partial x_i}$$

Handwritten example:  $\frac{df}{dx_1} = \frac{df}{da} \cdot \frac{da}{dx_1} = 1 \cdot 10 = 10$

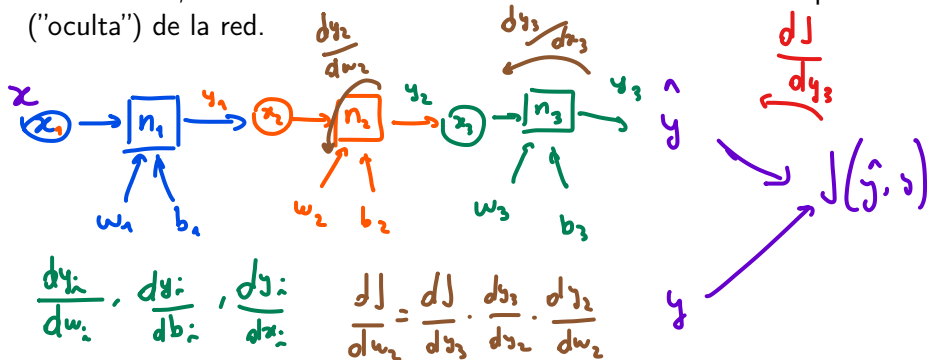
- La diferenciación automática se puede utilizar siempre que la función pueda representarse como un grafo de cómputo.
- La gran ganancia de este mecanismo está en que cada función sólo precisa saber cómo derivarse a sí misma, permitiendo OOP.

# 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:

*no lineal*  $\hat{y} = g(w \cdot x + b)$  *transf. afín = T.L. + despl.*  $x \rightarrow [n] \rightarrow y$

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  $z = w \cdot x + b$  y suponiendo conocido

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

$$\begin{aligned} \frac{\partial \hat{y}}{\partial z} &= g'(z) \in \mathbb{R} \\ \left[ \begin{aligned} \frac{\partial J}{\partial W} &= \frac{dJ}{d\hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial W} = dY \cdot g'(z) \cdot x \in \mathbb{R} \\ \frac{\partial J}{\partial b} &= \frac{dJ}{d\hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b} = dY \cdot g'(z) \cdot 1 \in \mathbb{R} \end{aligned} \right. \end{aligned}$$

## Redes neuronales (II)

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

$$\frac{\partial J}{\partial w_i} = dY \cdot g'(z) \cdot x_i$$

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

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \rightarrow \boxed{n} \rightarrow y$$

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

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

$$\begin{aligned} \text{si } W &\in \mathbb{R}^{1 \times n} \\ \Rightarrow \nabla_J(W) &\in \mathbb{R}^{1 \times n} \end{aligned}$$

$$\frac{\partial J}{\partial b} = dY \cdot g'(z) \in \mathbb{R}^1$$

# Rredes 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:

$$\hat{y} = g(\underbrace{W \cdot x + b}_{\in \mathbb{R}^k})$$



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  $b$ ?

$$\underbrace{\overbrace{W}^{k \times n} \cdot \overbrace{x}^{n \times 1}}_{k \times 1} + \underbrace{\overbrace{b}^{k \times 1}}_{k \times 1} = z = \begin{pmatrix} z_1 \\ \vdots \\ z_k \end{pmatrix} \xrightarrow{g} g(z) = \begin{pmatrix} g(z_1) \\ \vdots \\ g(z_k) \end{pmatrix}$$

# Redes neuronales (IV)

$$\frac{\partial J}{\partial b_i} = dY_i \cdot g'(z_i) \cdot 1$$

$$\begin{pmatrix} 0 \\ 1 \end{pmatrix} \odot \begin{pmatrix} c \\ d \end{pmatrix} = \begin{pmatrix} a \cdot c \\ b \cdot d \end{pmatrix}$$

En el caso de  $b$  es simple:

$$\begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} \quad \frac{\partial J}{\partial b} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial b} = \begin{pmatrix} dY_1 \\ \vdots \\ dY_k \end{pmatrix} \odot \begin{pmatrix} g'(z_1) \\ \vdots \\ g'(z_k) \end{pmatrix} \odot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = dY \odot g'(z)$$

Ahora para cada elemento de  $W$  tenemos:

$$\frac{\partial J}{\partial w_i} = dY_i \cdot g'(z_i) \cdot x^T$$

$$\begin{aligned} \frac{\partial J}{\partial W} &= \begin{pmatrix} \frac{\partial J}{\partial W_{1,1}} & \cdots & \frac{\partial J}{\partial W_{1,n}} \\ \frac{\partial J}{\partial W_{2,1}} & \cdots & \frac{\partial J}{\partial W_{2,n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial J}{\partial W_{k,1}} & \cdots & \frac{\partial J}{\partial W_{k,n}} \end{pmatrix} = \begin{pmatrix} -\nabla_J(W_{1,:}) \\ \vdots \\ -\nabla_J(W_{k,:}) \end{pmatrix} = \begin{bmatrix} dY_1 \cdot g'(z_1) \cdot x^T \\ \vdots \\ dY_k \cdot g'(z_k) \cdot x^T \end{bmatrix} = \\ &= \begin{pmatrix} dY_1 \\ \vdots \\ dY_k \end{pmatrix} \odot \begin{pmatrix} g'(z_1) \\ \vdots \\ g'(z_k) \end{pmatrix} \cdot x^T = \underbrace{dY \odot g'(z)}_{\text{green}} \cdot x^T \end{aligned}$$

# 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.

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

*Handwritten notes:*  
- Red arrows pointing from  $\frac{dJ}{d\hat{y}}$  to  $dY$  and from  $\frac{\partial J}{\partial y_k}$  to  $dY_k$ .  
- Red boxes around  $W_{1,j}$  and  $W_{k,j}$ .  
- Red box around  $z_j$ .  
- Red box around  $g'(z)$  in the second equation.

$$= \langle dY \odot g'(z), W_{:,j} \rangle = W_{j,:}^T \cdot dY \odot g'(z)$$

En forma vectorizada:

$$\langle a, b \rangle = a^T \cdot b = b^T \cdot a$$

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

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