

Análisis Matemático para Inteligencia Artificial

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

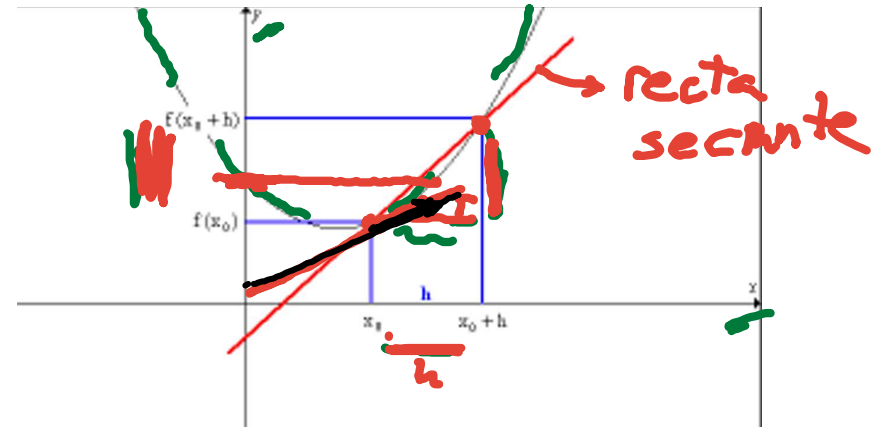
Especialización en Inteligencia Artificial

27/5/2022

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 $h \rightarrow 0$, ...

$$f: \mathbb{R} \rightarrow \mathbb{R} \\ x \mapsto y = f(x)$$



... esto nos conduce a la definición de derivada de f en x_0 :

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



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 $n = 1$,

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

Handwritten examples:

- $f(x) = x^2$ (with $n=1$)
- $f(x, y) = x^2 + y^2$ (with $n=2$)
- $f(t) = (t, 2t^2)$ (with $m=2, n=1$)
- $f(x, y) = (x^2, 0, y^2 - 1)$ (with $m=3, n=2$)

Level sets for $f(x, y) = x^2 + y^2$:

- $k=1: 1 = x^2 + y^2$
- $k=4: 4 = x^2 + y^2$

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.

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:

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

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

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:

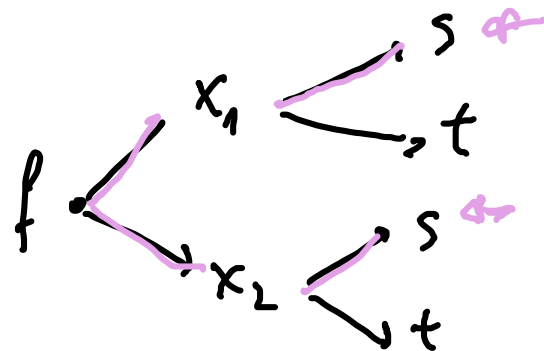
$$J_f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \rightarrow \begin{matrix} \nabla f_1 \\ \vdots \\ \nabla f_m \end{matrix} \quad J_f (m \times n)$$

Regla de la Cadena en forma matricial

Sea $f(\underbrace{x_1(s, t)}_{f_1}, \underbrace{x_2(s, t)}_{f_2})$

$$\frac{\partial f}{\partial s} = \left(\frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} \right) + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s}$$

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



Y luego

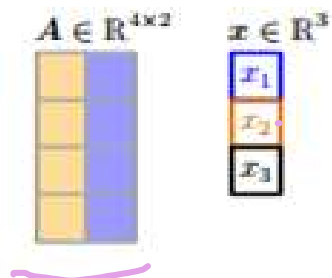
$$\frac{df}{d(s, t)} = \frac{df}{dx} \frac{dx}{d(s, t)} = \underbrace{\left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]}_{\nabla f(x_1, x_2)} \begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix}$$

Recordemos reglas de derivación:

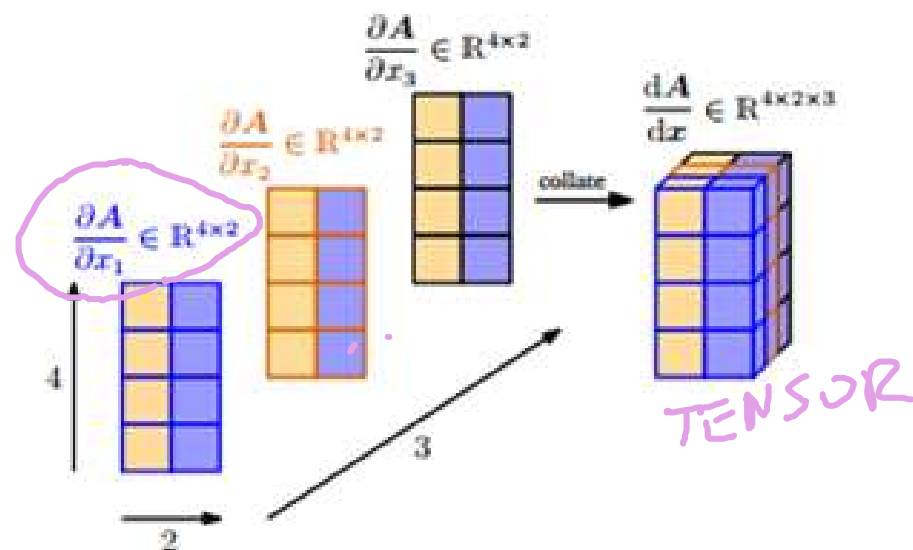
- $\frac{\partial (f+g)(s)}{\partial s} = \frac{\partial f}{\partial s} + \frac{\partial g}{\partial s}$
- $\frac{\partial (fg)(s)}{\partial s} = \frac{\partial f}{\partial s} g(s) + f(s) \frac{\partial g}{\partial s}$

Derivada de matrices

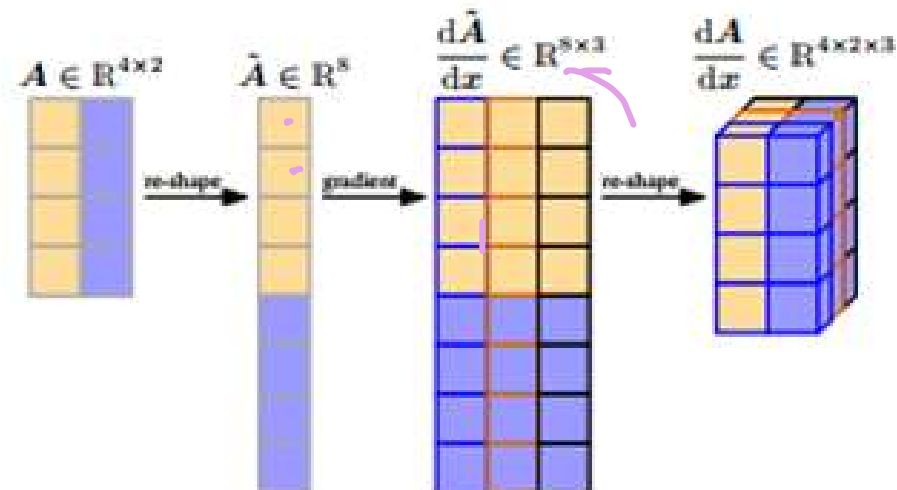
$A \times$



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 4×2 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:

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

$\frac{\partial^2 f}{\partial x_1^2}$ derivada segunda
 derivado resp a x_1 2 veces
 derivado resp. a x_n

$$\frac{\partial^2 f}{\partial x_1^2} = \frac{\partial}{\partial x_1} \left(\frac{\partial f}{\partial x_1} \right)$$

$$\frac{\partial^2 f}{\partial x_1 \partial x_n} = \frac{\partial}{\partial x_1} \left(\frac{\partial f}{\partial x_n} \right)$$

Ejemplo $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

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

• Calculo $\frac{\partial f}{\partial x}(x, y) = 2y^2$

$$\frac{\partial f}{\partial y}(x, y) = 4xy$$

$$\frac{\partial^2 f}{\partial x \partial y}(x, y) = 4y$$

$$\frac{\partial^2 f}{\partial x^2}(x, y) = 0$$

$$\frac{\partial^2 f}{\partial y \partial x}(x, y) = 4y$$

$$\frac{\partial^2 f}{\partial y^2}(x, y) = 4x$$

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

$H = H^T$ simétrica

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

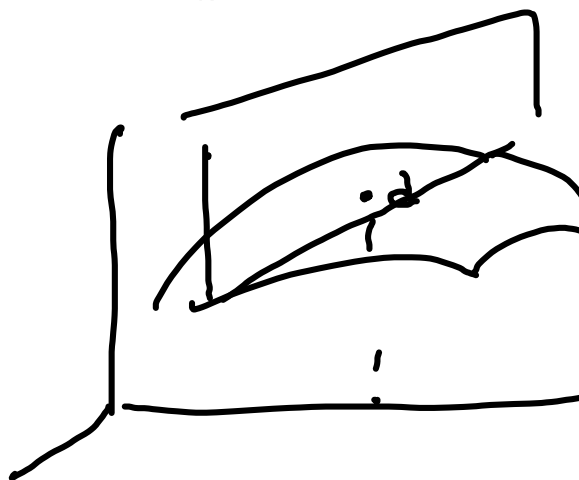
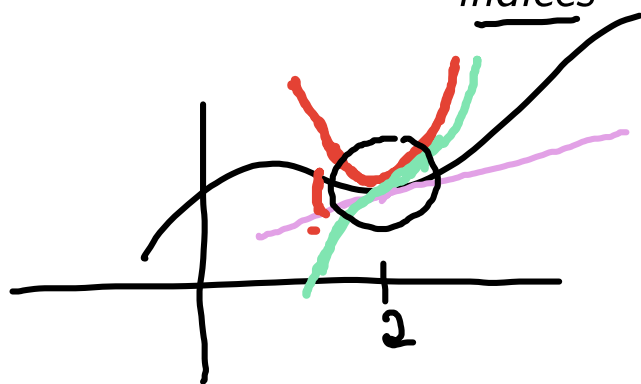
Relacionado a F. Cuadráticas

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 :

$$a = (a_1, \dots, a_n)$$

$$P_k(x) = f(a) + \underbrace{\sum_{i=1}^n \frac{\partial f}{\partial x_i}(a)(x_i - a_i)}_{\text{aprox. 1º ord.}} + \frac{1}{2!} \sum_{i,j=1}^n \frac{\partial^2 f}{\partial x_i \partial x_j}(a) \underbrace{(x_i - a_i)(x_j - a_j)}_{\text{aprox. 2º ord.}} + \dots + \frac{1}{3!} \sum \partial^3 + \dots + \frac{1}{k!} \sum_{\text{indices}} \frac{\partial^k f}{\partial x_{i_1} \dots \partial x_{i_k}}(a)(x_{i_1} - a_{i_1}) \dots (x_{i_k} - a_{i_k})$$

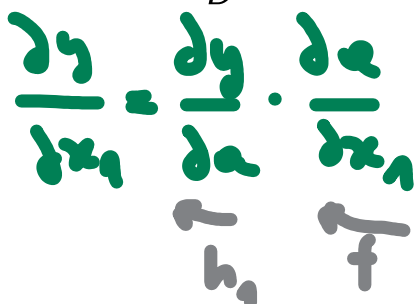


Diferenciación Automática

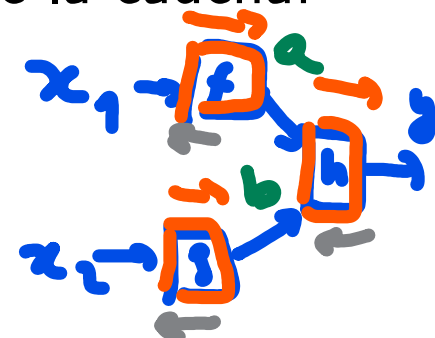
Sean, para una función f :

- 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:

$$\frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial a} \cdot \frac{\partial a}{\partial x_1}$$


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



- 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 Cube:
    def forward(x):
        self.last_x = x
        return x**3
    def backwards():
        return (self.last_x**2) * 3
```

Diferenciación Automática: ejemplo

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

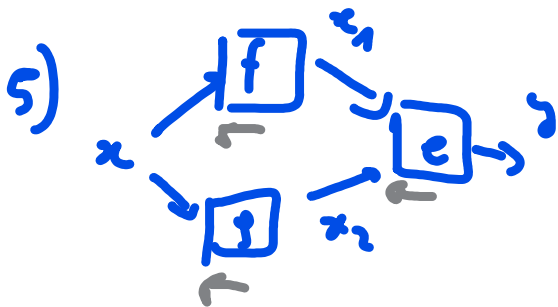
1) $x \rightarrow [f] \rightarrow z$ $\frac{dz}{dx} = f'(x)$

3) $x_1, x_2 \rightarrow [e] \rightarrow z$ $\nabla_e(x_1, x_2) = \left(\frac{\partial e}{\partial x_1}, \frac{\partial e}{\partial x_2} \right)$
 $\hookrightarrow \frac{\partial e}{\partial x_1} = \nabla_{e_1}(x_1, x_2)$

2) $x \rightarrow [g] \rightarrow z$ $\frac{dz}{dx} = g'(x)$

4) $x \rightarrow [f] \rightarrow y \rightarrow [g] \rightarrow z$ $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = g'(y) \cdot f'(x)$
 $y = g \circ f(x) = (3x)^2$

requiere forward de f



$$\frac{dz}{dx} = \frac{dz}{dx_1} \cdot \frac{dx_1}{dx} + \frac{dz}{dx_2} \cdot \frac{dx_2}{dx}$$

$$= \nabla_{e_1}(x_1, x_2) \cdot f'(x) + \nabla_{e_2}(x_1, x_2) \cdot g'(x)$$

$y = (3x) \cdot (x^2)$

$$\frac{dy}{dx} = \frac{dy}{df} \cdot \frac{df}{dx} + \frac{dy}{dg} \cdot \frac{dg}{dx}$$

Backpropagation

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

$$y = f(x, w, b)$$

¿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.

