

# Análisis Matemático para Inteligencia Artificial

Martín Errázquin (merrazquin@fi.uba.ar)

Especialización en Inteligencia Artificial

Clase 5

# Análisis Matemático

- 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

$x_0$ :

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

# 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**,  $n > 1$ .
- Si  $m > 1$  diremos que es una función
  - **vectorial**, si  $n = 1$ ,
  - **campo vectorial**,  $n > 1$ .

**Conjuntos de Nivel** Dada  $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}$  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:

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

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

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:

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

# Matriz Hessiana

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

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

Una aplicación común del Hessiano es el polinomio de Taylor de orden 2 para campos escalares. 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 orden 2:

$$P_2(x) = f(a) + \nabla f(a)^T (x - a) + \frac{1}{2} (x - a)^T H_f(a) (x - a)$$

# Regla de la Cadena en forma matricial

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

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

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

Y luego

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

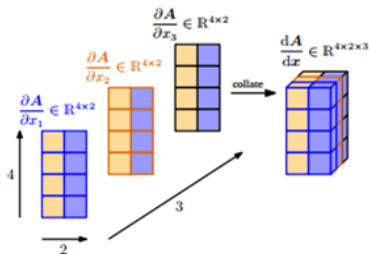
Recordemos reglas de derivación:

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

# Derivada de matrices

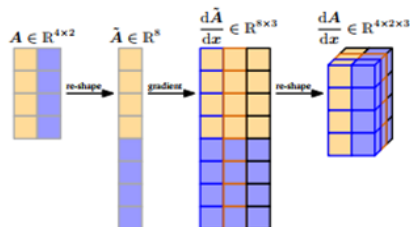
$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$


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.

$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$

(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:

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

- 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.

# Diferenciación automática: idea gráfica

# 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: Resumen

Sea  $X \in \mathbb{R}^n$ ,  $W \in \mathbb{R}^{k \times n}$ ,  $b \in \mathbb{R}^k$  y  $g : \mathbb{R} \rightarrow \mathbb{R}$  no lineal, si se recibe  $\frac{dJ}{dy} \in \mathbb{R}^k$  entonces:

$$\frac{dJ}{db} = dY \odot g'(z)$$

$$\frac{dJ}{dW} = dY \odot g'(z) \cdot X^T$$

$$\frac{dJ}{dX} = W^T \cdot dY \odot g'(z)$$

donde  $z = W \cdot X + b \in \mathbb{R}^k$  e  $y = g(z)$  con  $g$  aplicada elemento a elemento.

Luego como  $y^{(m)} = x^{(m+1)}$ , entonces  $\frac{dJ}{dX^{(m+1)}} = \frac{dJ}{dy^{(m)}}$  que es lo que le pasamos a la capa anterior.

# Redes neuronales: Bosquejo de código

```
class Layer:
    ...
    def forward(self, X):
        self.last_x = X
        self.last_z = self.W @ X + self.b
        self.last_y = self.g.f(self.last_z)
        return self.last_y

    def backwards(self, dY):
        db = dY * self.g.df(self.last_z)
        dW = db @ self.last_x.T
        dX = self.W.T @ db

        self.W, self.b = self.optimizer.step(self.W, self.b, dW, db)

        return dX
    ...
```

Para no ocupar demasiado tiempo de clase con la derivación del desarrollo que lleva a las fórmulas que vimos antes, dejamos una playlist con videos de youtube en el campus en la sección de esta semana.

¡Aprovecharla!