

Bibliografía

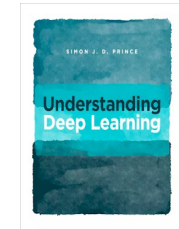
Tema 2 – Optimización y Regularización (Parte 3)

Aprendizaje Automático II - Grado en Inteligencia Artificial
Universidad Rey Juan Carlos

Iván Ramírez Díaz
ivan.ramirez@urjc.es

José Miguel Buenaposada Biencinto
josemiguel.buenaposada@urjc.es

- **Understanding Deep Learning**. Capítulo 7.



- **Deep Learning for Computer Vision**. Lecture 4. Stanford University. 2017. Curso en youtube.
- **Deep Learning: CS 182 2021**. Lecture 5. Sergey Levine. UC Berkeley. Curso en youtube.

Algoritmo de descenso de gradiente

- **Paso 0**. Inicializar los parámetros Φ_0 .
- **Repetir**:
 - **Paso 1**. Calcular derivadas de **la función de coste** con respecto a los parámetros Φ

$$\left. \frac{\partial J(\Phi)}{\partial \Phi} \right|_{\Phi=\Phi_t} = \begin{bmatrix} \frac{\partial J(\Phi)}{\partial \Phi_0} \\ \frac{\partial J(\Phi)}{\partial \Phi_1} \\ \vdots \\ \frac{\partial J(\Phi)}{\partial \Phi_k} \end{bmatrix}$$

- **Paso 2**. Actualizar los parámetros de acuerdo con:

$$\Phi_{t+1} \leftarrow \Phi_t - \alpha \left. \frac{\partial J(\Phi)}{\partial \Phi} \right|_{\Phi=\Phi_t}$$

donde el escalar positivo α determina la magnitud del cambio.

Implementación del descenso de gradiente

- **Problema 1**. Calcular derivadas de **la función de coste** con respecto a los parámetros

$$\left. \frac{\partial J(\Phi)}{\partial \Phi} \right|_{\Phi=\Phi_t} = \begin{bmatrix} \frac{\partial J(\Phi)}{\partial \Phi_0} \\ \frac{\partial J(\Phi)}{\partial \Phi_1} \\ \vdots \\ \frac{\partial J(\Phi)}{\partial \Phi_k} \end{bmatrix}$$

¿Por qué es un problema?

- Muchos modelos son muy complicados con multitud de parámetros:

$$y' = \phi'_0 + \phi'_1 a [\psi_{10} + \psi_{11} a [\theta_{10} + \theta_{11} x] + \psi_{12} a [\theta_{20} + \theta_{21} x] + \psi_{13} a [\theta_{30} + \theta_{31} x]] \\ + \phi'_2 a [\psi_{20} + \psi_{21} a [\theta_{10} + \theta_{11} x] + \psi_{22} a [\theta_{20} + \theta_{21} x] + \psi_{23} a [\theta_{30} + \theta_{31} x]] \\ + \phi'_3 a [\psi_{30} + \psi_{31} a [\theta_{10} + \theta_{11} x] + \psi_{32} a [\theta_{20} + \theta_{21} x] + \psi_{33} a [\theta_{30} + \theta_{31} x]]$$

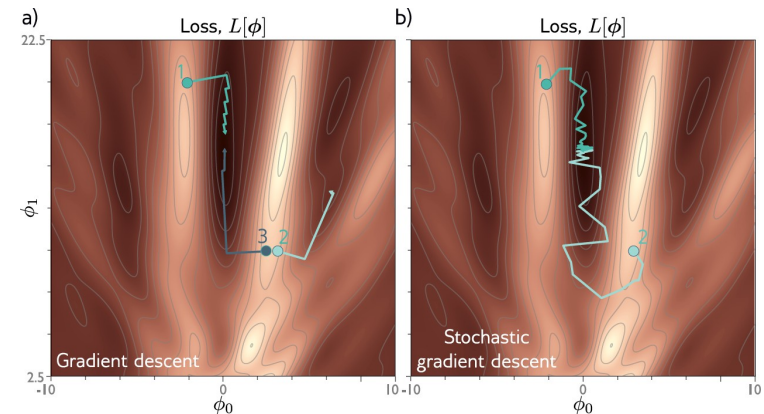
- Hay que calcular la derivada:
 - Para cada parámetro
 - Para cada muestra en el mini-batch
 - Para cada iteración del SGD

2.4 Cálculo eficiente del gradiente

- Conceptos matemáticos
- Retropropagación del gradiente (Backpropagation)
- Diferenciación algorítmica

Implementación del descenso de gradiente

- **Problema 2.** Inicialización



¿Dónde deberíamos inicializar los parámetros antes de arrancar el SGD?

2.4 Cálculo eficiente del gradiente

- **Conceptos matemáticos**
- Retropropagación del gradiente (Backpropagation)
- Diferenciación algorítmica

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = \frac{dg(x)}{dx} \cdot \left. \frac{df(y)}{dy} \right|_{y=g(x)}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = 4 \cdot \left. \frac{df(y)}{dy} \right|_{y=g(x)}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = 4 \cdot (2 \cdot y)|_{y=g(x)}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = 4 \cdot (2 \cdot y)|_{y=1+4x}$$

Recordatorio: regla de la cadena

- Tenemos una composición de funciones sobre escalares:

$$x \xrightarrow{g} y \xrightarrow{f} z$$

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

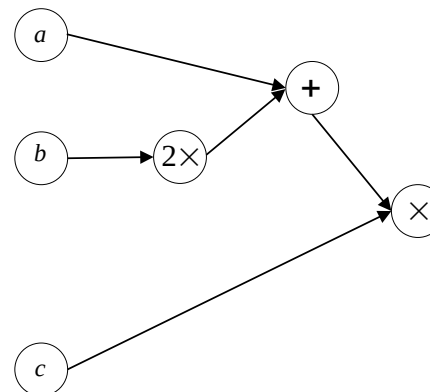
Ejemplo: $g(x) = 1 + 4x$; $f(y) = y^2$

$$\frac{d}{dx}(f(g(x))) = 4 \cdot 2 \cdot (1 + 4x) = 8 + 32x$$

Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a, b, c) = (a + 2 \cdot b) \cdot c$$



Grafos de cómputo

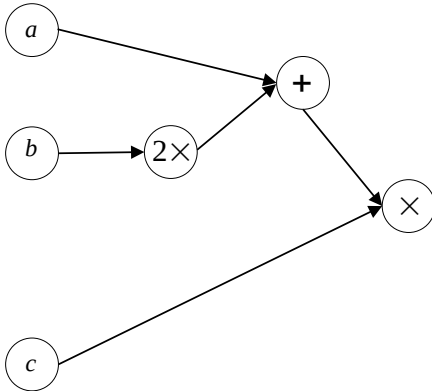
¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$



Grafos de cómputo

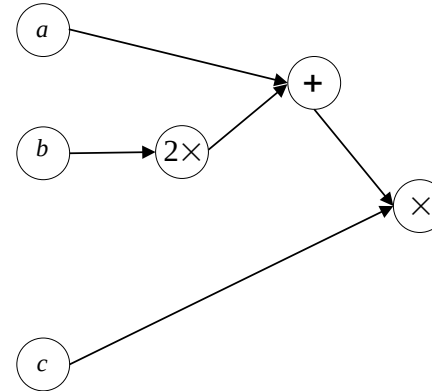
¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$



Calculemos: $\frac{\partial J}{\partial a}$

Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

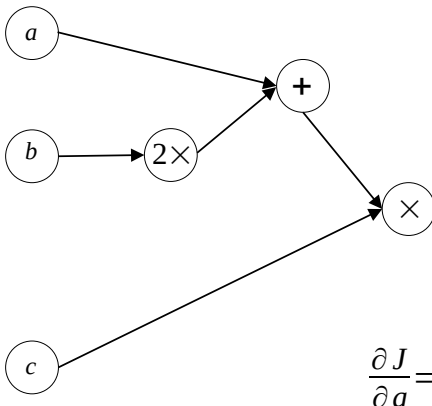
Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:



$$\frac{\partial J}{\partial a} =$$

Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

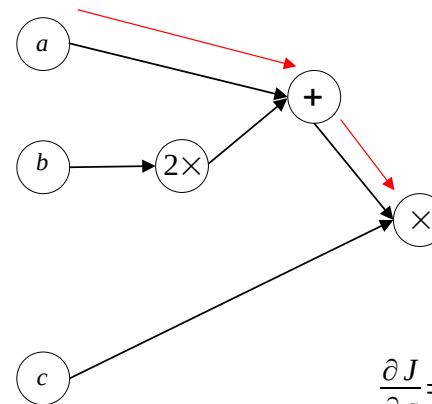
Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:



$$\frac{\partial J}{\partial a} = \frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x} + \dots$$

Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

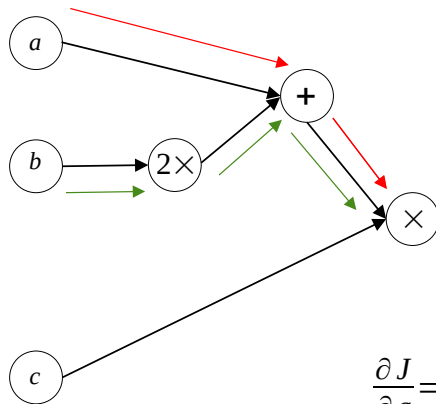
$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = \boxed{\frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x}} + \boxed{\frac{\partial b}{\partial a} \cdot \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x}} + \dots$$



Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

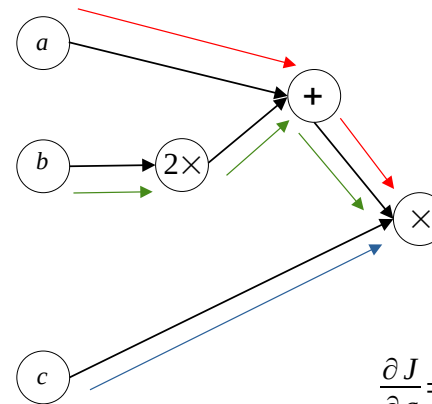
$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = \boxed{\frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x}} + \boxed{\frac{\partial b}{\partial a} \cdot \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x}} + \boxed{\frac{\partial c}{\partial a} \cdot \frac{\partial f}{\partial c}}$$



Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

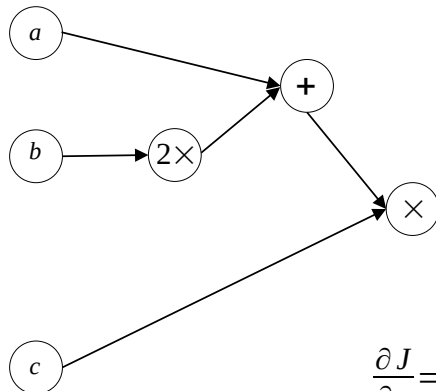
$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = \boxed{\frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x}} + \cancel{\boxed{\frac{\partial b}{\partial a} \cdot \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x}}} + \cancel{\boxed{\frac{\partial c}{\partial a} \cdot \frac{\partial f}{\partial c}}}$$

0 0



Grafos de cómputo

¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

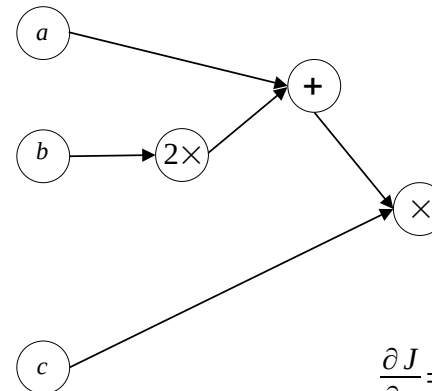
$$J(a,b,c)=f(g(a,h(b)),c)$$

$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

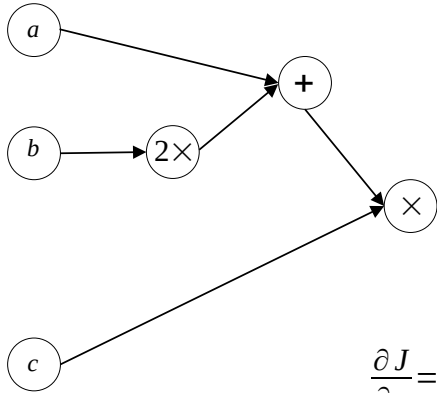
Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = \boxed{\frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x}}$$



Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

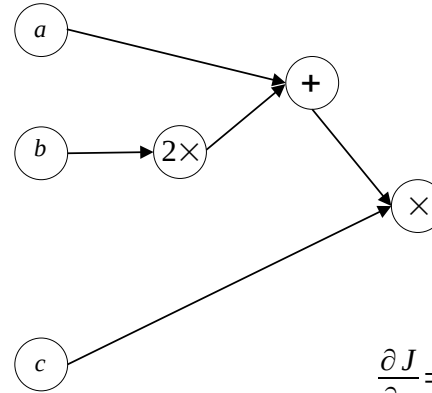
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = \frac{\partial g}{\partial a} \Big|_{s=h(b)} \cdot \frac{\partial f}{\partial x} \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

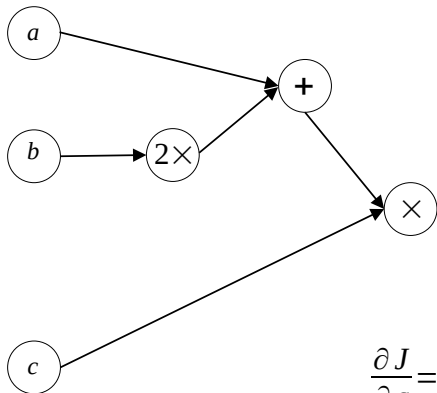
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = 1 \cdot \frac{\partial f}{\partial x} \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

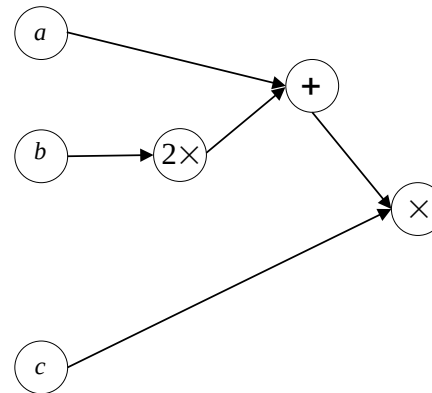
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial a}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial a} = 1 \cdot c$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

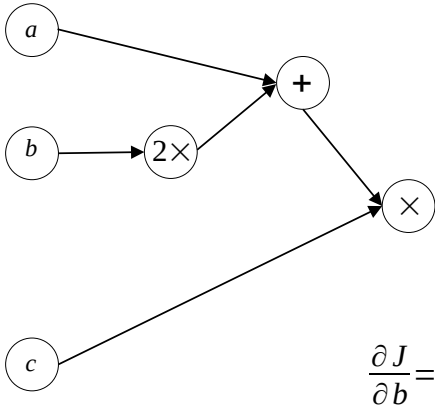
Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

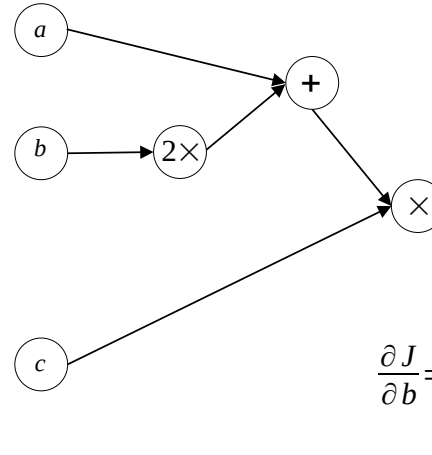
$$f(x,c)=x \cdot c \quad g(a,s)=a+s \quad h(b)=2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} =$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

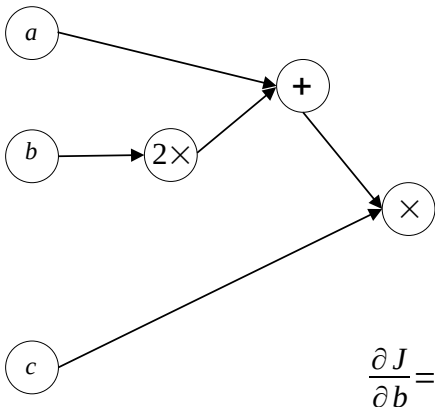
$$f(x,c)=x\cdot c \quad g(a,s)=a+s \quad h(b)=2\cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = \underbrace{\frac{\partial a}{\partial b}}_0 \cdot \frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x} + \underbrace{\frac{\partial c}{\partial b}}_0 \cdot \frac{\partial f}{\partial c}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

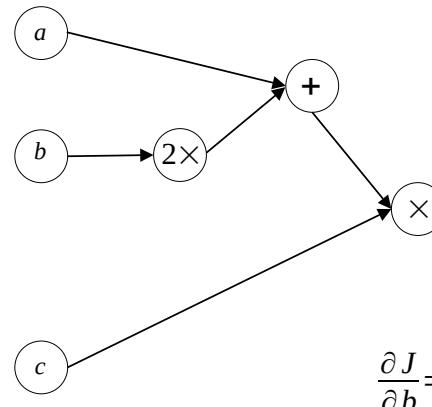
$$f(x,c)=x \cdot c \quad g(a,s)=a+s \quad h(b)=2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c)=(a+2\cdot b)\cdot c$$

Composición de funciones:

$$J(a,b,c)=f(g(a,h(b)),c)$$

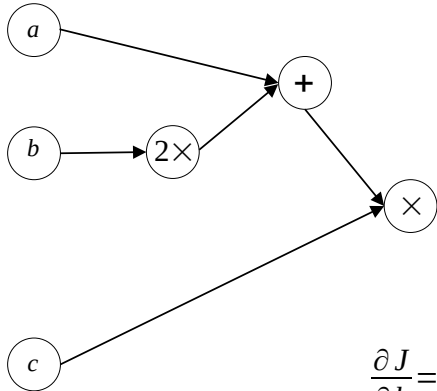
$$f(x,c)=x \cdot c \quad g(a,s)=a+s \quad h(b)=2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \bigg|_{s=h(b)} \cdot \frac{\partial f}{\partial x} \bigg|_{x=g(a, h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

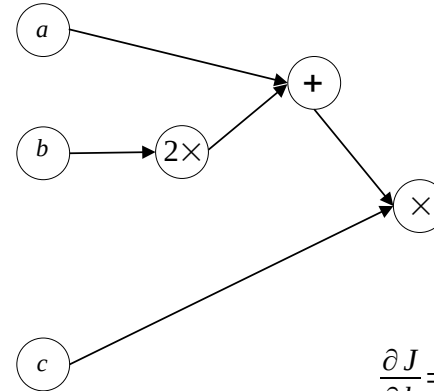
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = 2 \cdot \frac{\partial g}{\partial s} \Big|_{s=h(b)} \cdot \frac{\partial f}{\partial x} \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

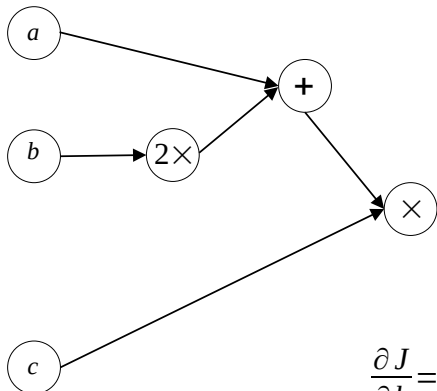
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = 2 \cdot 1 \cdot \frac{\partial f}{\partial x} \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

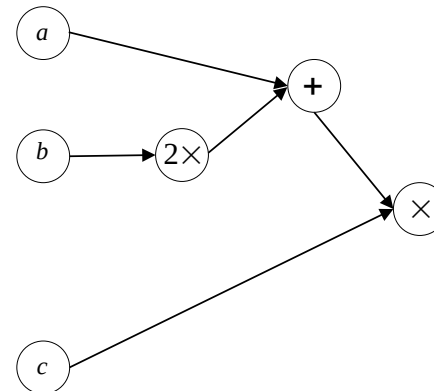
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial b}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial b} = 2 \cdot 1 \cdot c = 2 \cdot c$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

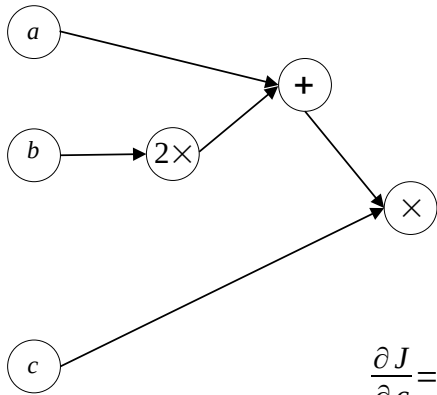
Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial c}$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

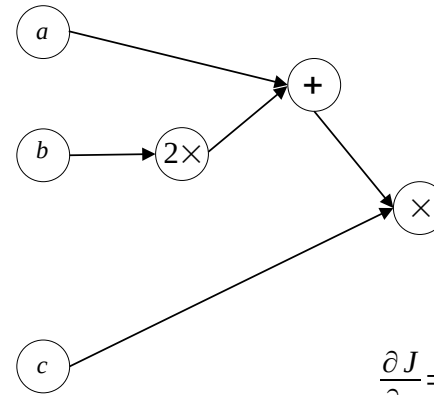
Calculemos: $\frac{\partial J}{\partial c}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial c} = \frac{\partial a}{\partial c} \cdot \frac{\partial g}{\partial a} \cdot \frac{\partial f}{\partial x} + \frac{\partial b}{\partial c} \cdot \frac{\partial h}{\partial b} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial f}{\partial x} + \frac{\partial f}{\partial c}$$

0 0

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

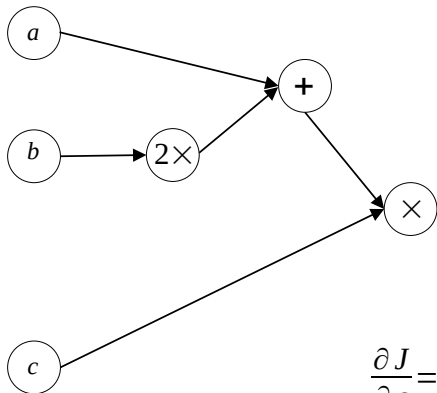
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial c}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial c} = \frac{\partial f}{\partial c} \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

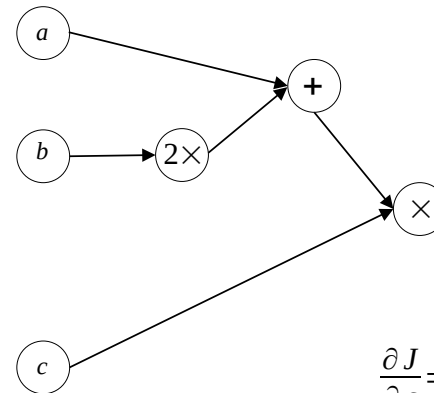
$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial c}$

La derivada usa todos los caminos:

$$\frac{\partial J}{\partial c} = x \Big|_{x=g(a,h(b))}$$

Grafos de cómputo



¿qué **expresión** calcula el grafo?

$$J(a,b,c) = (a+2 \cdot b) \cdot c$$

Composición de funciones:

$$J(a,b,c) = f(g(a,h(b)),c)$$

$$f(x,c) = x \cdot c \quad g(a,s) = a+s \quad h(b) = 2 \cdot b$$

Calculemos: $\frac{\partial J}{\partial c}$

La derivada usa todos los caminos:

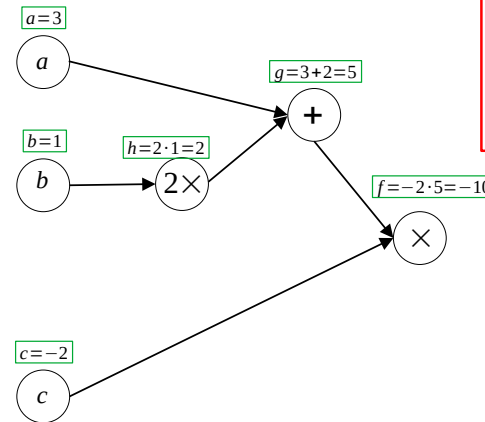
$$\frac{\partial J}{\partial c} = a+2b$$

2.4 Cálculo eficiente del gradiente

- Conceptos matemáticos
- Retropropagación del gradiente (Backpropagation)
- Diferenciación algorítmica

Retropropagación en grafos de cómputo

Fase 1: ejecución del grafo (forward)



$$J(a, b, c) = (a + 2 \cdot b) \cdot c$$

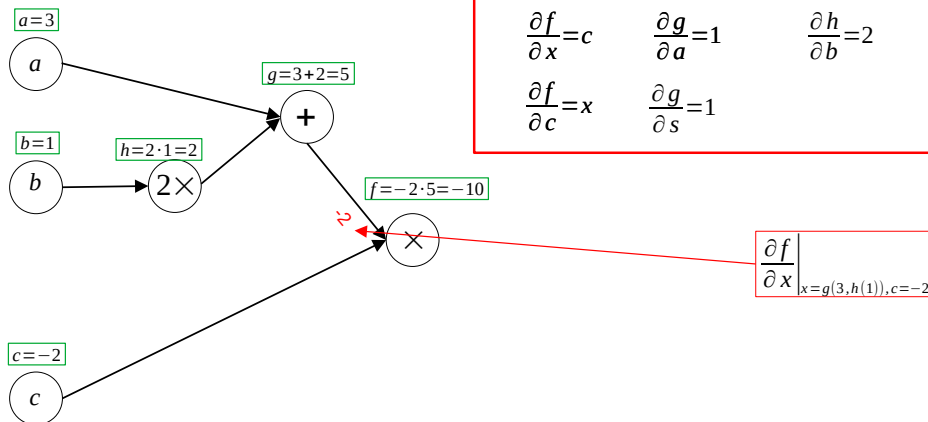
$$f(x, c) = x \cdot c \quad g(a, s) = a + s \quad h(b) = 2 \cdot b$$

$$\frac{\partial f}{\partial x} = c \quad \frac{\partial g}{\partial a} = 1 \quad \frac{\partial h}{\partial b} = 2$$

$$\frac{\partial f}{\partial c} = x \quad \frac{\partial g}{\partial s} = 1$$

Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



$$J(a, b, c) = (a + 2 \cdot b) \cdot c$$

$$f(x, c) = x \cdot c \quad g(a, s) = a + s \quad h(b) = 2 \cdot b$$

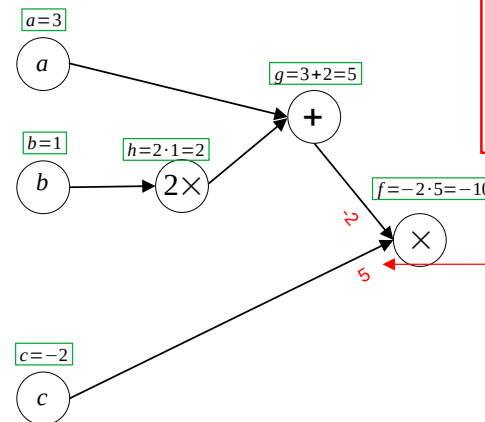
$$\frac{\partial f}{\partial x} = c \quad \frac{\partial g}{\partial a} = 1 \quad \frac{\partial h}{\partial b} = 2$$

$$\frac{\partial f}{\partial c} = x \quad \frac{\partial g}{\partial s} = 1$$

$$\left. \frac{\partial f}{\partial x} \right|_{x=g(3, h(1)), c=-2}$$

Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



$$J(a, b, c) = (a + 2 \cdot b) \cdot c$$

$$f(x, c) = x \cdot c \quad g(a, s) = a + s \quad h(b) = 2 \cdot b$$

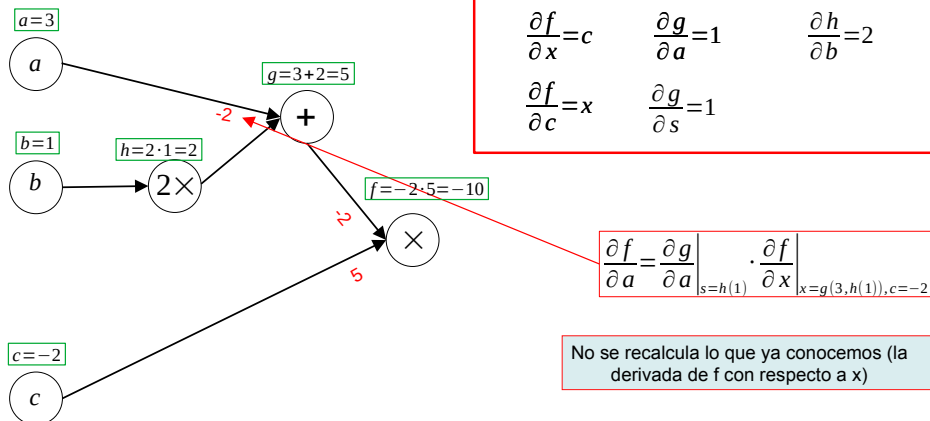
$$\frac{\partial f}{\partial x} = c \quad \frac{\partial g}{\partial a} = 1 \quad \frac{\partial h}{\partial b} = 2$$

$$\frac{\partial f}{\partial c} = x \quad \frac{\partial g}{\partial s} = 1$$

$$\left. \frac{\partial f}{\partial c} \right|_{x=g(3, h(1)), c=-2}$$

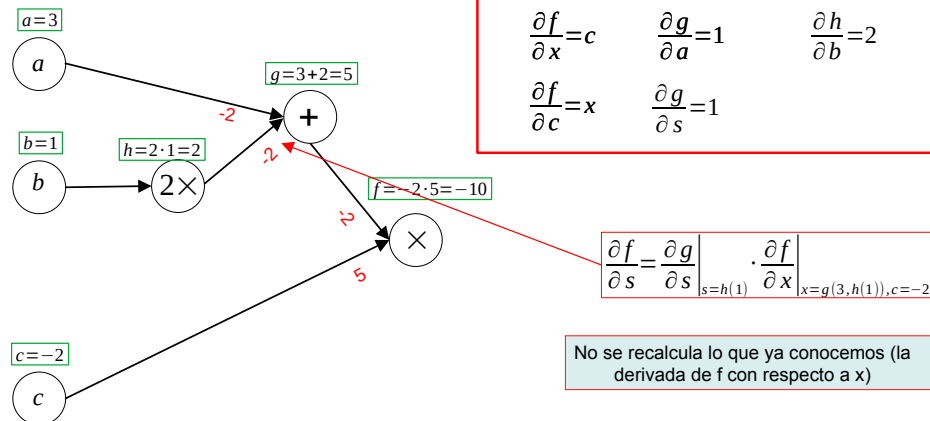
Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



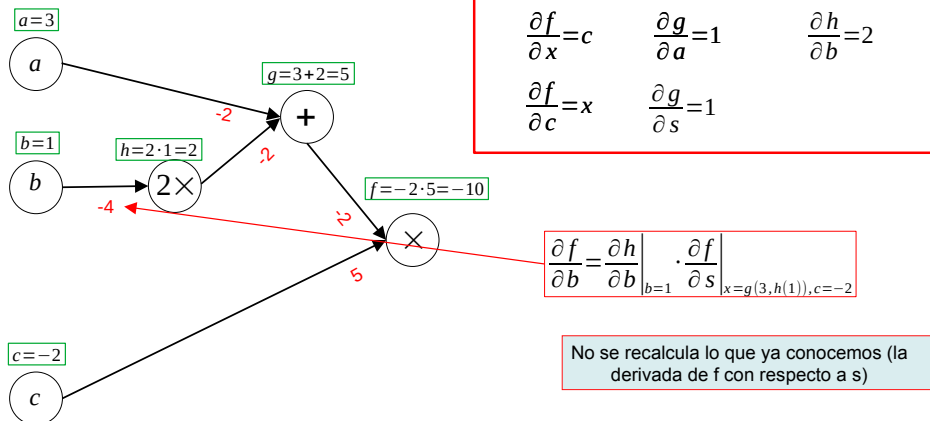
Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



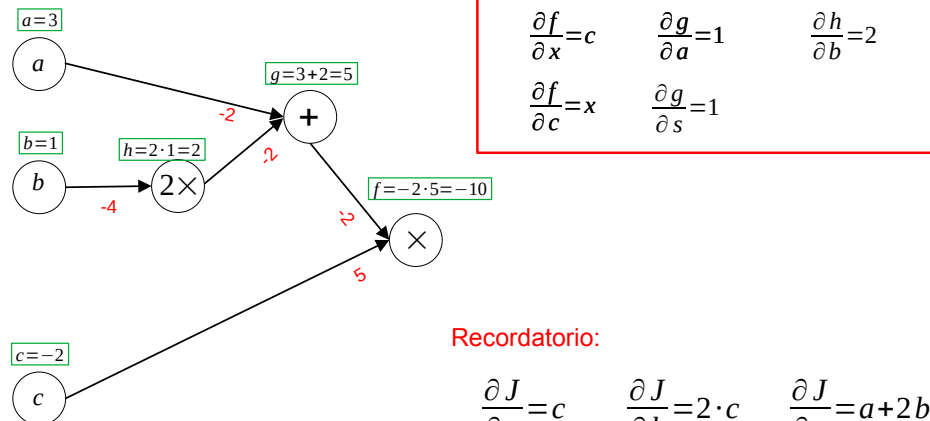
Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



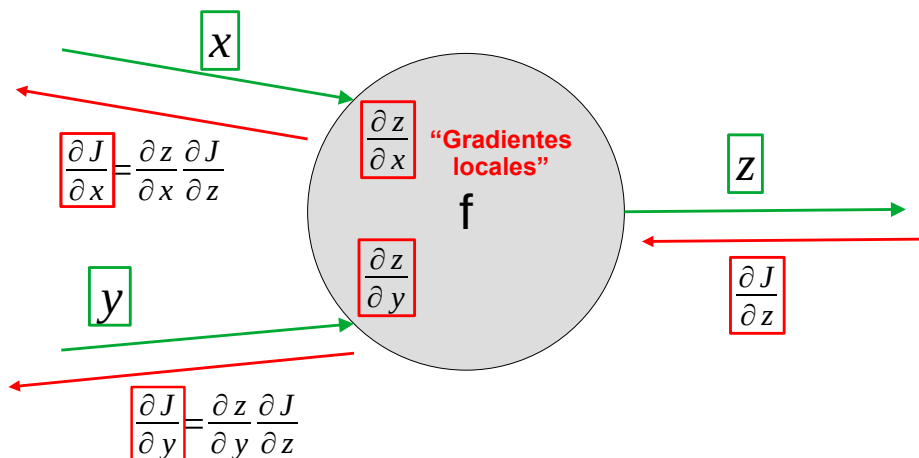
Retropropagación en grafos de cómputo

Fase 2: retropropagación (backward)



Retropropagación en grafos de cómputo

- En cada nodo únicamente necesitamos calcular las derivadas locales con respecto a sus entradas:



Algoritmo de retropropagación

- Forward:** Se ejecuta el grafo de cómputo dadas unas entradas
- Backward:**

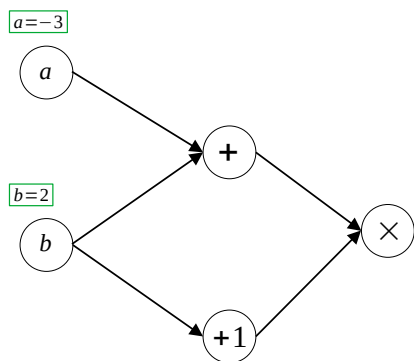
- Inicializar $\delta = \frac{dJ}{dz^{(n)}}$ Salida de la última operación del grafo de cómputo

- Para cada f con entrada x_f desde el final al principio:

$$\delta \leftarrow \frac{df}{dx_f} \cdot \delta$$

Grafos de cómputo: Ejemplo 2

¿qué **expresión** calcula el grafo?



$$J(a, b) = (a+b) \cdot (b+1)$$

$$g = a+b$$

$$h = b+1$$

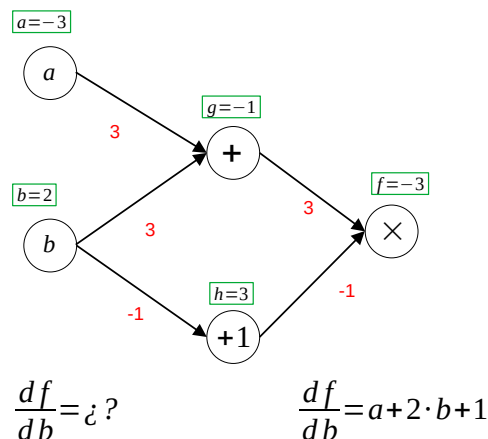
$$f = g \cdot h$$

$$\frac{\partial f}{\partial g} = h \quad \frac{\partial f}{\partial h} = g$$

$$\frac{\partial h}{\partial b} = 1 \quad \frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} = 1$$

Grafos de cómputo: Ejemplo 2

¿qué **expresión** calcula el grafo?



$$J(a, b) = (a+b) \cdot (b+1)$$

$$g = a+b$$

$$h = b+1$$

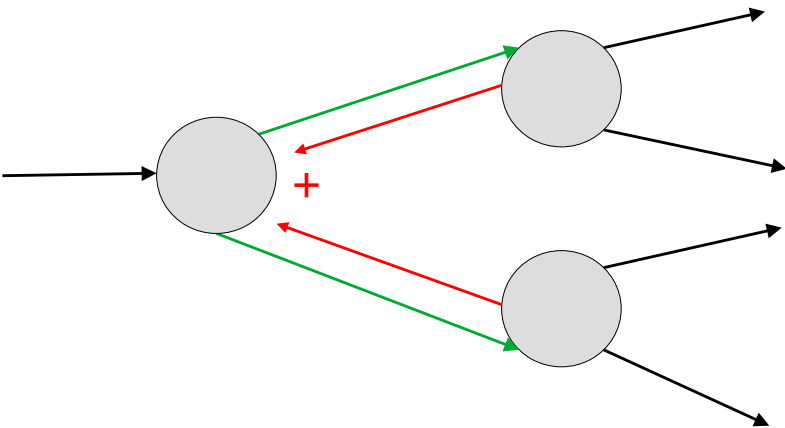
$$f = g \cdot h$$

$$\frac{\partial f}{\partial g} = h \quad \frac{\partial f}{\partial h} = g$$

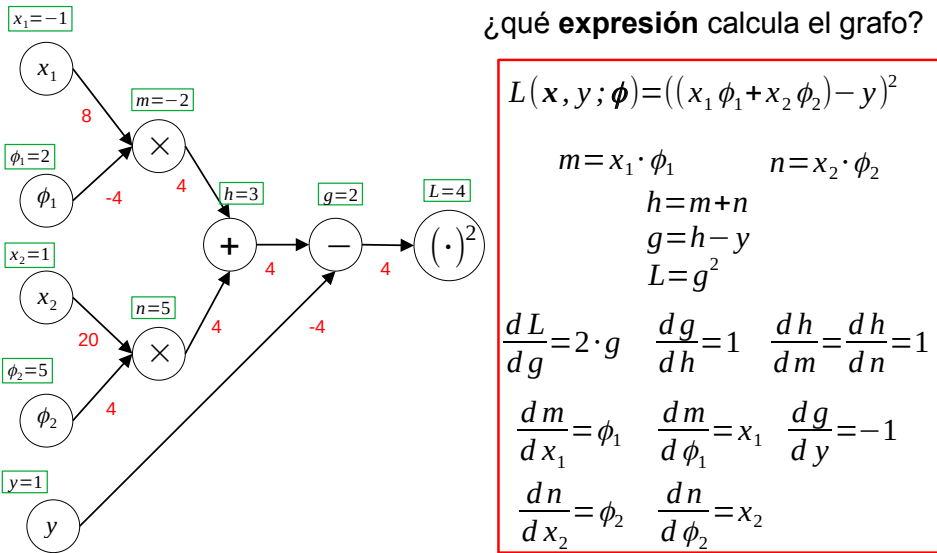
$$\frac{\partial h}{\partial b} = 1 \quad \frac{\partial g}{\partial a} = \frac{\partial g}{\partial b} = 1$$

Retropropagación en grafos de cómputo

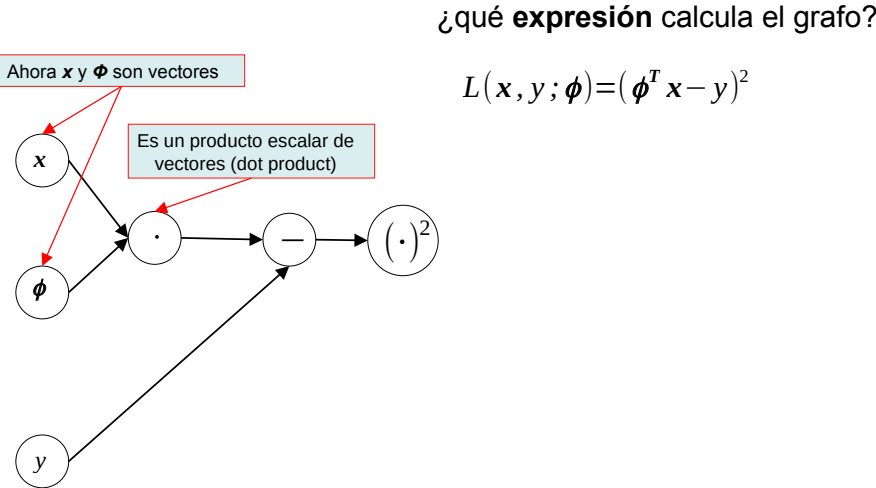
- Los gradientes se suman en las ramificaciones:



Grafos de cómputo: Ejemplo 3

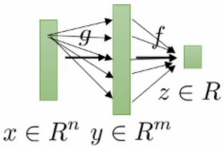


Grafos de cómputo: vectores y matrices



Recordatorio: regla de la cadena

- Si tenemos una composición de funciones **sobre vectores**:



- Por la regla de la cadena de Cálculo:

$$\frac{dz}{dx} (f(g(x))) = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

Jacobiano de g Jacobiano de f

$y \in R^m$ $\frac{dz}{dy} \in R^m$ $\frac{dy}{dx} \in R^{n \times m}$

Recordatorio: regla de la cadena

- Por la regla de la cadena de Cálculo:

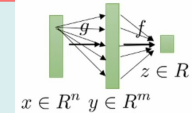
$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Diagram illustrating the chain rule for the derivative of a composition of functions. The diagram shows three vertical bars representing the domains: $y \in R^m$, $\frac{dz}{dy} \in R^m$, and $\frac{dy}{dx} \in R^{n \times m}$. Red arrows point from the labels "Jacobiano de g" and "Jacobiano de f" to the corresponding terms in the equation.

Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Diagram illustrating the chain rule for the derivative of a composition of functions. The diagram shows three vertical bars representing the domains: $y \in R^m$, $\frac{dz}{dy} \in R^m$, and $\frac{dy}{dx} \in R^{n \times m}$. Red arrows point from the labels "Jacobiano de g" and "Jacobiano de f" to the corresponding terms in the equation.

Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} w_{11}x_1 + w_{12}x_2 + w_{13}x_3 \\ w_{21}x_1 + w_{22}x_2 + w_{23}x_3 \end{bmatrix}$$

Diagram illustrating the composition of functions g and f. g maps x in R^n to y in R^m, and f maps y in R^m to z in R. The diagram shows a network of nodes and connections representing the functions.

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Diagram illustrating the chain rule for the derivative of a composition of functions. The diagram shows three vertical bars representing the domains: $y \in R^m$, $\frac{dz}{dy} \in R^m$, and $\frac{dy}{dx} \in R^{n \times m}$. Red arrows point from the labels "Jacobiano de g" and "Jacobiano de f" to the corresponding terms in the equation.

Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} w_{11}x_1 + w_{12}x_2 + w_{13}x_3 \\ w_{21}x_1 + w_{22}x_2 + w_{23}x_3 \end{bmatrix}$$

Diagram illustrating the composition of functions g and f. g maps x in R^n to y in R^m, and f maps y in R^m to z in R. The diagram shows a network of nodes and connections representing the functions.

- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

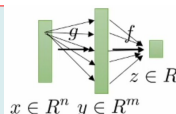
Diagram illustrating the chain rule for the derivative of a composition of functions. The diagram shows three vertical bars representing the domains: $y \in R^m$, $\frac{dz}{dy} \in R^m$, and $\frac{dy}{dx} \in R^{n \times m}$. Red arrows point from the labels "Jacobiano de g" and "Jacobiano de f" to the corresponding terms in the equation.

Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

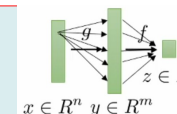
Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} \frac{\partial}{\partial x_1}(f(g(x))) \\ \frac{\partial}{\partial x_2}(f(g(x))) \\ \frac{\partial}{\partial x_3}(f(g(x))) \end{bmatrix} = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

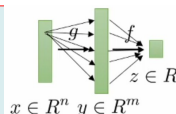
Recordatorio: regla de la cadena

Ejemplo:

$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

Recordatorio: regla de la cadena

Ejemplo:

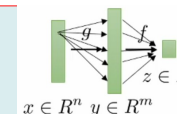
$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$

Regla de la cadena:

$$\frac{d}{dx}(f(g(x))) = \frac{g(x)}{dx} \cdot \frac{df(y)}{dy} \Big|_{y=g(x)}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

Recordatorio: regla de la cadena

Ejemplo:

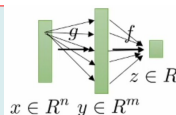
$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$

Regla de la cadena:

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \cdot \frac{df(y)}{dy} \Big|_{y=g(x)}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

Recordatorio: regla de la cadena

Ejemplo:

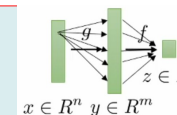
$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$

Regla de la cadena:

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \cdot \frac{df(y)}{dy} \Big|_{y=g(x)}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

Recordatorio: regla de la cadena

Ejemplo:

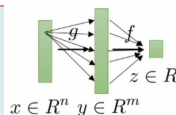
$$g(x) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad f(y) = \begin{bmatrix} 4 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$f(g(x)) = 4 \cdot (w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3) + 3 \cdot (w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3)$$

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} 4 \cdot w_{11} + 3 \cdot w_{21} \\ 4 \cdot w_{12} + 3 \cdot w_{22} \\ 4 \cdot w_{13} + 3 \cdot w_{23} \end{bmatrix}$$

Regla de la cadena:

$$\frac{d}{dx}(f(g(x))) = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \cdot \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$



- Por la regla de la cadena de Cálculo:

$$\frac{d}{dx}(f(g(x))) = \frac{dz}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

Jacobiano de g

Jacobiano de f

$$y \in R^m \quad \frac{dz}{dy} \in R^m \quad \frac{dy}{dx} \in R^{n \times m}$$

Derivadas con matrices

Function escalar (1D) $f[\mathbf{a}]$ de un vector \mathbf{a}

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}$$

$$\frac{\partial f}{\partial \mathbf{a}} = \begin{bmatrix} \frac{\partial f}{\partial a_1} \\ \frac{\partial f}{\partial a_2} \\ \frac{\partial f}{\partial a_3} \\ \frac{\partial f}{\partial a_4} \end{bmatrix}$$

Derivadas con matrices

Function escalar (1D) $f[]$ de un a matriz \mathbf{A}

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{bmatrix} \quad \frac{\partial f}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \frac{\partial f}{\partial a_{12}} & \frac{\partial f}{\partial a_{13}} \\ \frac{\partial f}{\partial a_{21}} & \frac{\partial f}{\partial a_{22}} & \frac{\partial f}{\partial a_{23}} \\ \frac{\partial f}{\partial a_{31}} & \frac{\partial f}{\partial a_{32}} & \frac{\partial f}{\partial a_{33}} \\ \frac{\partial f}{\partial a_{41}} & \frac{\partial f}{\partial a_{42}} & \frac{\partial f}{\partial a_{43}} \end{bmatrix}$$

Derivadas con matrices

Function vectorial $f[]$ de un vector \mathbf{a}

$$\mathbf{f} = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} \quad \frac{\partial \mathbf{f}}{\partial \mathbf{a}} = \begin{bmatrix} \frac{\partial f_1}{\partial a_1} & \frac{\partial f_2}{\partial a_1} & \frac{\partial f_3}{\partial a_1} \\ \frac{\partial f_1}{\partial a_2} & \frac{\partial f_2}{\partial a_2} & \frac{\partial f_3}{\partial a_2} \\ \frac{\partial f_1}{\partial a_3} & \frac{\partial f_2}{\partial a_3} & \frac{\partial f_3}{\partial a_3} \\ \frac{\partial f_1}{\partial a_4} & \frac{\partial f_2}{\partial a_4} & \frac{\partial f_3}{\partial a_4} \end{bmatrix}$$

Vectores vs Matrices

Derivadas escalares:

$$f_3 = \beta_3 + \omega_3 h_3 \quad \frac{\partial f_3}{\partial h_3} = \frac{\partial}{\partial h_3} (\beta_3 + \omega_3 h_3) = \omega_3$$

Derivadas de matrices

$$\mathbf{f}_3 = \boldsymbol{\beta}_3 + \boldsymbol{\Omega}_3 \mathbf{h}_3 \quad \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} = \frac{\partial}{\partial \mathbf{h}_3} (\boldsymbol{\beta}_3 + \boldsymbol{\Omega}_3 \mathbf{h}_3) = \boldsymbol{\Omega}_3^T$$

Vectores vs Matrices

Derivadas escalares:

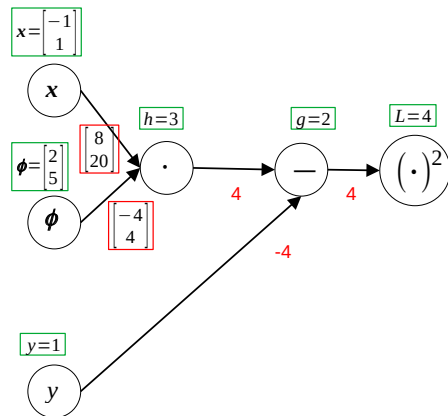
$$f_3 = \beta_3 + \omega_3 h_3 \quad \frac{\partial f_3}{\partial \beta_3} = \frac{\partial}{\partial \omega_3} \beta_3 + \omega_3 h_3 = 1$$

Derivadas de matrices

$$\mathbf{f}_3 = \boldsymbol{\beta}_3 + \boldsymbol{\Omega}_3 \mathbf{h}_3 \quad \frac{\partial \mathbf{f}_3}{\partial \boldsymbol{\beta}_3} = \frac{\partial}{\partial \beta_3} (\boldsymbol{\beta}_3 + \boldsymbol{\Omega}_3 \mathbf{h}_3) = \mathbf{I}$$

Grafos de cómputo

¿qué **expresión** calcula el grafo?



$$L(x, y; \phi) = (\phi^T x - y)^2$$

$$h = \phi^T x$$

$$g = h - y$$

$$L = g^2$$

$$\frac{dL}{dg} = 2 \cdot g \quad \frac{dg}{dh} = 1 \quad \frac{dg}{dy} = -1$$

$$\frac{dh}{dx} = \phi \quad \frac{dh}{d\phi} = x$$

Diferenciación algorítmica

- Los frameworks de deep learning calculan derivadas automáticamente
- Únicamente se especifica el modelo y la función de pérdida
- ¿Cómo? **Diferenciación Algorítmica**
 - Cada componente sabe como calcular su propia derivada
 - Una función lineal “sabe” cómo calcular su derivada de la salida con respecto a su entrada
 - Una función lineal “sabe” cómo calcular su derivada de la salida con respecto a sus parámetros
 - Se especifica el orden de los componentes (grafo de cómputo)
 - Puede calcular la cadena de derivadas

2.4 Cálculo eficiente del gradiente

- Conceptos matemáticos
- Retropropagación del gradiente (Backpropagation)
- Diferenciación algorítmica**

Algoritmo de retropropagación

- Forward:** Se ejecuta el grafo de cómputo dadas unas entradas
- Backward:**

– Inicializar $\delta = \frac{dJ}{dz^{(n)}}$ Salida de la última operación del grafo de cómputo

– Para cada f con entrada x_f y parámetros ϕ_f desde el final al principio:

Derivada de J con respecto a los parámetros de este nodo

$$\frac{df}{d\phi_f} = \frac{df}{dz_f} \cdot \delta$$

$$\delta \leftarrow \frac{df}{dx_f} \cdot \delta$$