

2) Back Propagation

1

$L(\theta)$ θ parameter

L función de error que depende de los parámetros

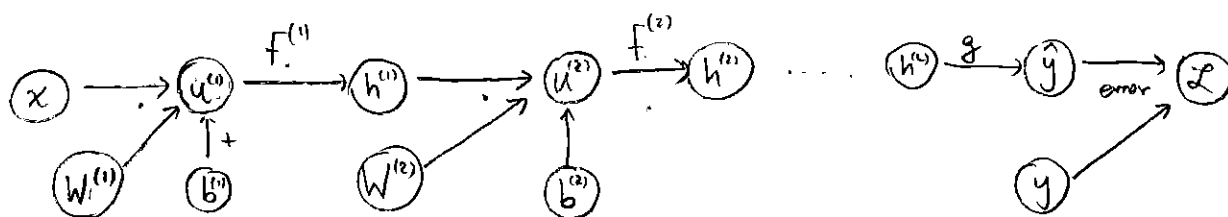
$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\} \quad \hat{y}^{(i)} = \text{Forward}(x^{(i)})$$

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \text{error}(\hat{y}^{(i)}, y^{(i)})$$

queremos $\nabla_{\theta} L(\theta) = \left(\frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_r} \right)$ y evaluarlo en un punto $\hat{\theta}$ para poder derivar

Back propagation es simplemente un método para calcular $\nabla_{\theta} L(\theta)$ y evaluarlo eficientemente en el computador.

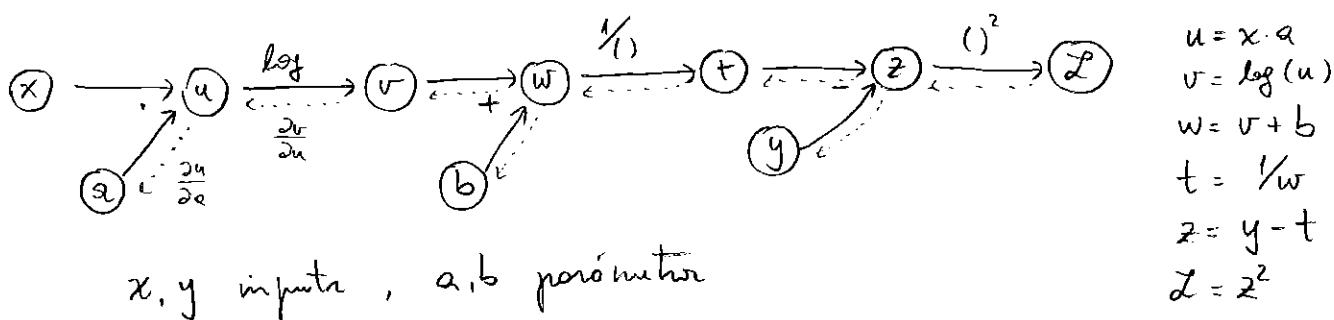
Primer concepto importante: Computational graph



Forma de organizar toda las operaciones que engloban lo por lo hacia adelante y el cálculo de la pérdida.

¿cómo calcular $\frac{\partial L}{\partial W_{ij}^{(2)}}$? $L = \text{error}(g(f^{(2)}(f^{(1)}(f^{(1-1)}(f^{(2)}(h^{(1)}W^{(2)} + b^{(2)})W^{(3)}))))$
no parece una tarea fácil :c

Idea: usar aproximadamente la regla de la cadena!
consideremos un computational graph simple (solo escalar, sin matrices)



¿ $\frac{\partial \mathcal{L}}{\partial a}$? idea 1) desarrollar la definición de \mathcal{L}

$$\mathcal{L} = \left(y - \frac{1}{b + \log(x \cdot a)} \right)^2$$

idea 2) usar la regla de la cadena "desde atrás" recursivamente

¿ $\frac{\partial \mathcal{L}}{\partial a}$? $= \frac{\partial \mathcal{L}}{\partial u} \cdot \frac{\partial u}{\partial a} = \frac{\partial \mathcal{L}}{\partial u} \cdot x$

si tuviera $\frac{\partial \mathcal{L}}{\partial u}$ computarlo podría
computar $\frac{\partial \mathcal{L}}{\partial a}$ fácilmente

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial v}{\partial u} = \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{1}{u}$$

$$\frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial w} \cdot \frac{\partial w}{\partial v} = \frac{\partial \mathcal{L}}{\partial w} \cdot 1$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial t} \cdot \frac{\partial t}{\partial w} = \frac{\partial \mathcal{L}}{\partial t} \cdot \frac{-1}{w^2}$$

$$\frac{\partial \mathcal{L}}{\partial t} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial t} = \frac{\partial \mathcal{L}}{\partial z} \cdot -1$$

$$\frac{\partial \mathcal{L}}{\partial z} = 2z \quad \text{¡into!} \checkmark$$

¿ $\frac{\partial \mathcal{L}}{\partial b}$? MUY FÁCIL!

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial w} \cdot \frac{\partial w}{\partial b} = \frac{\partial \mathcal{L}}{\partial w} \cdot 1$$

¡yo lo tengo
pre calculado!!

Idea general: para cada nodo n en el grafo

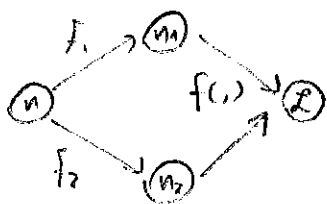


primero calculo $\frac{\partial \mathcal{L}}{\partial n'}$ y luego lo uso para computar

$$\frac{\partial \mathcal{L}}{\partial n} = \frac{\partial \mathcal{L}}{\partial n'} \cdot \frac{\partial n'}{\partial n}$$

esto define un uso de cada
número de calcular y de valores
dados los datos computados durante
la pasada forward

¿ Qué pasa si tengo algo como esto?



$$n_1 = f_1(n)$$

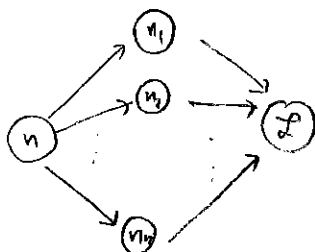
$$n_2 = f_2(n)$$

$$\mathcal{L} = f(n_1, n_2)$$

¿ $\frac{\partial \mathcal{L}}{\partial n}$?

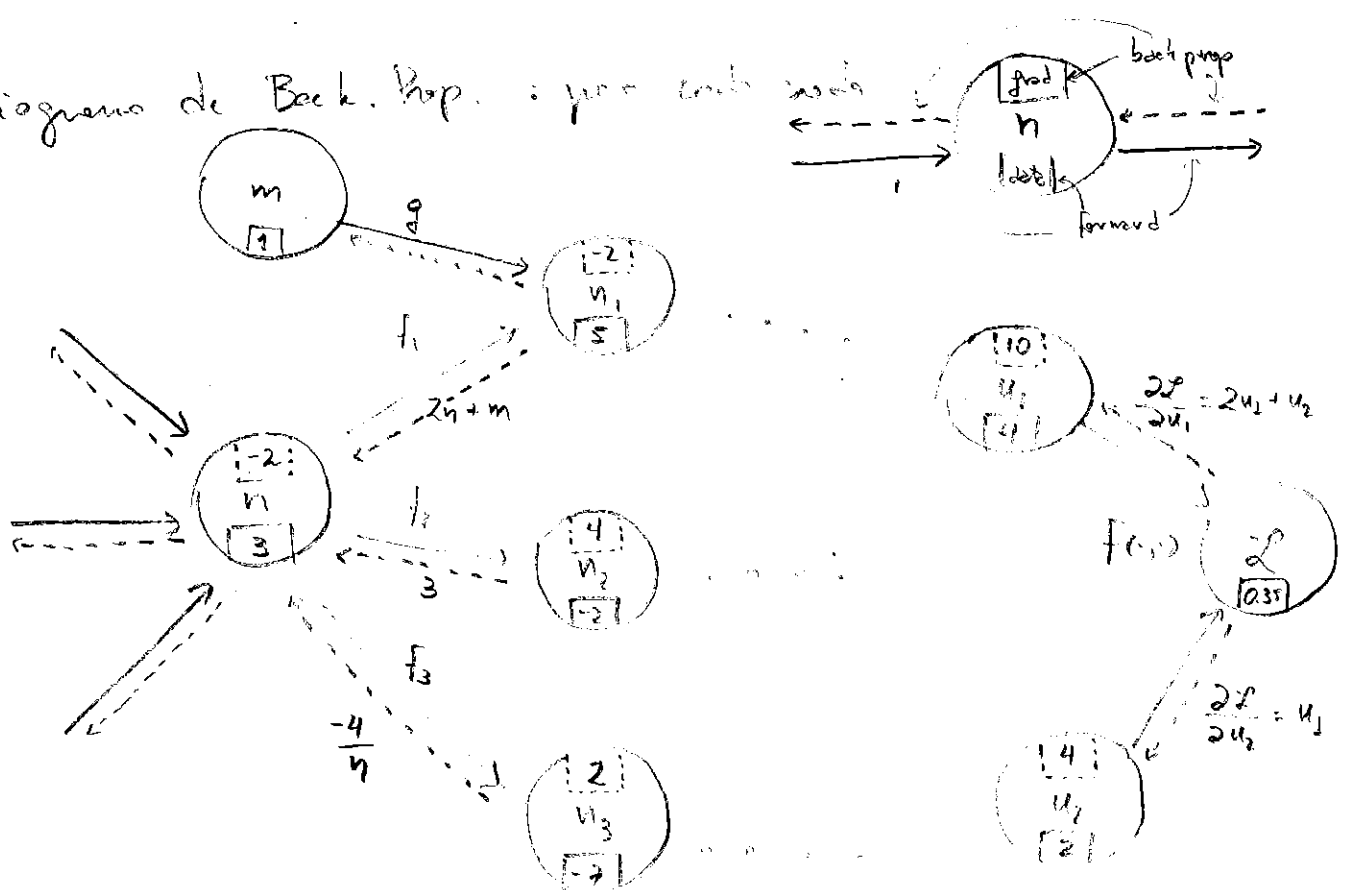
$$\frac{\partial \mathcal{L}}{\partial n} = \frac{\partial \mathcal{L}}{\partial n_1} \cdot \frac{\partial n_1}{\partial n} + \frac{\partial \mathcal{L}}{\partial n_2} \cdot \frac{\partial n_2}{\partial n}$$

y generalize



$$\frac{\partial \mathcal{L}}{\partial n} = \sum_{i=1}^k \frac{\partial \mathcal{L}}{\partial n_i} \cdot \frac{\partial n_i}{\partial n}$$

Diagrama de Back. Pop. : por cada individuo ←



parameters are "keys" of sets (not like in actualized words)

Algorithme de backpropagation (récursif)

Naciter :

- Variables:
- G : computational graph (DAG), - 2 variable objectives
 - $data[n]$: para cada nodo $n \in G$ almacena el valor computado en forward (incluyendo inputs, parámetros y nodos intermedios)

Colours :

- grad [n] para cada parámetro

$$\text{grad}[\mathcal{L}] := 1$$

$$\text{rdprop}(n):$$

- ```

if (n) :
 if good[n] != NULL return good[n]
 else good[n] = NULL

```

$\text{gcd}[n] = \text{MML}$   
 con  $n_1, n_2, \dots, n_n$  los que dependen de  $n$  en  $G$

$$\text{grad}[n] := \sum_{i=1}^k \text{rbprop}(n_i) \cdot \left( \frac{\partial n}{\partial n_i} \right) (\text{data})$$

para cada parámetro  $p$ , ejecutar  $rbprop(p)$

gitarre in der Garage  
Tanz / Party / Theater  
Pyrotechnik / Laser  
automatisches

" $\frac{\partial n}{\partial n_i}$ " es la única parte no trivial del anterior código  
 y es lo que ayuda a hacer los librerías (TF, Theano, Pytorch)

## Algoritmo Back Propagation (iterativo)

Neuró: -

-  $G, L, \text{data}$  (igual que antes)

Cálculo:

-  $\text{grad}[n]$  para cada  $n$  en  $G$

bprop:

$\text{grad}[L] := 1$ ,  $\text{grad}[n] := 0$  para  $n \neq L$

mientras  $G$  tenga nodos sin marcar:

$n :=$  nodo en  $G$  con todas sus salidas salientes marcadas

para cada  $(n', n) \in G$ :

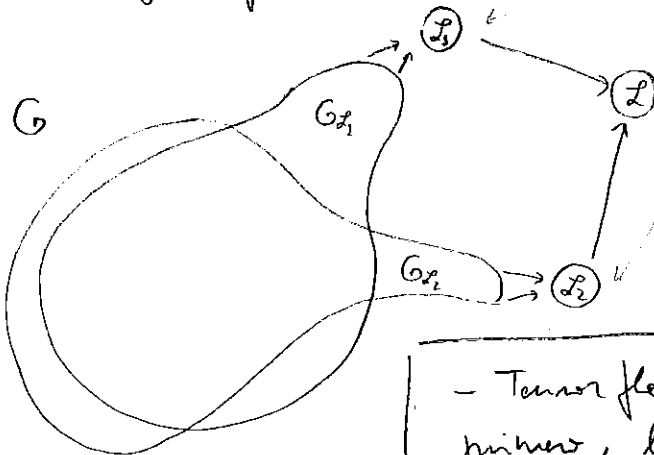
$\text{grad}[n'] += \text{grad}[n] \cdot \frac{\partial n}{\partial n'}(\text{data})$

marca  $(n', n)$

marca  $n$

aunque ningún  
 los únicos no trivial

Back Propagation puede pensarse más general



donde funciones sigmoideas (pérdida)  
 podemos hacer bprop sobre  $G_1$   
 y luego sobre  $G_2$ , incluso  
 podemos acumular los gradientes

Back Propagation es lo práctico:

- Tensorflow: se define el grafo de computación  
 primero, luego se "compila" lo que calcule los  
 derivados parciales y se almacenan en otros grafos

y luego se hacen bprop. - Pytorch: el grafo está implícito en las operaciones,  
 se realizan desde datos/parámetros, los derivados se calculan on the fly cuando  
 se hacen el bprop. (mucho más versátil!!!)