

Back-propagation Through Time (BPTT)

①

y_d → Value real (label/target)

L_t → loss function

\hat{y}_t → output

softmax neurons...

output

w_y → peso hidden-output

w_h → hidden

w_x → peso input-hidden

Input

Equations

$$L_T = \sum_t L_t$$

$$L_t = -y_t \log \hat{y}_t$$

$$\hat{y}_t = \text{softmax}(o_t) = \frac{e^{o_t}}{\sum_k e^{o_k}}$$

$$o_t = h_t \cdot w_y + b_y$$

$$h_t = \text{Tanh}[x_t \cdot w_x + h_{t-1} \cdot w_h + b_h]$$

$$L_T = \sum_t L_t$$

bias hidden

L_t

L_t

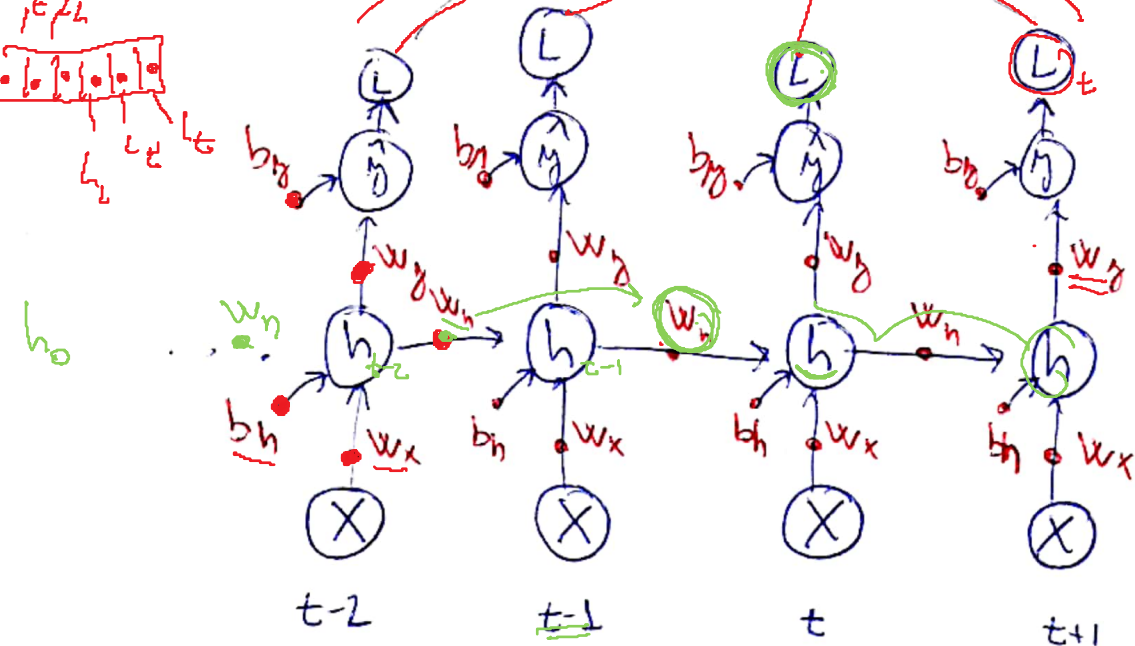
L_t

L_t

L_t

L_t

L_t



Parámetros a optimizar.

- 1º W_y → Peso hidden - output
- 2º b_y → bias output
- 3º W_h → Peso hidden - hidden
- 4º W_x → Peso input - hidden
- 5º b_h → bias hidden

¿Como varia mi L_T respecto de ellos?

$$\begin{aligned}
 \frac{\partial L_T}{\partial W_y} &= \frac{\partial}{\partial W_y} \left[\sum_t L_t \right] = \sum_t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial W_y}
 \end{aligned}$$

Los pesos no varían con los Δt
 sumo todos los t para tener una L global
 derivada de Loss function respecto de una softmax
 Ya fue demostrado $\left(\hat{y}_t - y_t \right)$

de ecuacion (4) en hoja (I)

$$o_t = W_y \cdot h_t + b_y$$

$$\frac{\partial o_t}{\partial W_y} = h_t$$

entonces llegamos a que

$$\frac{\partial L_t}{\partial W_y} = \left(\hat{y}_t - y_t \right) h_t \Rightarrow$$

$$\frac{\partial L_T}{\partial W_y} = \sum_t \left(\hat{y}_t - y_t \right) h_t \quad (6)$$

III

2º $\frac{\partial L_T}{\partial b_y}$ se calcula igual q 1º reemplazando la última $\frac{\partial o_t}{\partial w_y}$ por $\frac{\partial o_t}{\partial b_y}$

$$\frac{\partial L_T}{\partial b_y} = \sum_t (\hat{y}_t - y_t) \cdot 1 \quad \oplus$$

$$\frac{\partial o_t}{\partial b_y}$$

$$o_t = w_y \cdot h_t + b_y$$

3º $\frac{\partial L_T}{\partial w_h}$ ¿Cómo varia L_T cuando varía w_h ?

sol. $h_t = \text{Tanh} \begin{pmatrix} x w_x \\ h w_h \\ b_h \end{pmatrix}$

$$o_{t+1} = w_h h_{t+1} + b_y$$

$$\frac{\partial L_{t+1}}{\partial w_h} = \frac{\partial L_{t+1}}{\partial \hat{y}_{t+1}} \cdot \frac{\partial \hat{y}_{t+1}}{\partial o_{t+1}} \cdot \frac{\partial o_{t+1}}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial w_h} \quad \text{a resolver...}$$

Ya lo vieron

$$\left(\hat{y}_{t+1} - y_{t+1} \right)$$

de ecuación 4 en hoja I = w_y

$$h_{t+1} = \text{Tanh} (X_L w_x + h_t w_h + b_h)$$

$$\frac{\partial h_{t+1}}{\partial w_h} = 0$$

h_{t+1} es función de w_h y h_t

q es función de w_h y h_{t-2}

$$\frac{\partial h_{t+1}}{\partial w_h} = \frac{\partial h_{t+1}}{\partial h_t} \frac{\partial h_t}{\partial w_h} \Rightarrow \text{Podemos escribir:}$$

$$\frac{\partial h_{t+1}}{\partial w_h} = \sum_{k=1}^{t+1} \frac{\partial h_{t+1}}{\partial h_k} \frac{\partial h_k}{\partial w_h}$$

→ veamos que da eso

$$\left(\frac{\partial h_3}{\partial w_n} \right) = \sum_{k=1}^3 \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial w_n}$$

$$= \underbrace{\frac{\partial h_3}{\partial h_1} \frac{\partial h_1}{\partial w_n}}_{\substack{h_3 \text{ varia porque} \\ h_1 \text{ varia porque} \\ w_n \text{ varia}}} + \underbrace{\frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial w_n}}_{\substack{h_3 \text{ varia xq} \\ h_2 \text{ varia xq} \\ w_n \text{ varia}}} + \underbrace{\frac{\partial h_3}{\partial h_3} \frac{\partial h_3}{\partial w_n}}_{\substack{h_3 \text{ varia xq} \\ w_n \text{ varia}}}$$

(IV)

$$h_3 = x_3 \cdot w_n + w_h \cdot h_2 + b_h$$

Recapitulamos ecuacion (8) de nota (III)

$$\frac{\partial L_{t+1}}{\partial w_n} = \frac{\partial L_{t+1}}{\partial \hat{y}_{t+1}} \cdot \frac{\partial \hat{y}_{t+1}}{\partial o_{t+1}} \cdot \frac{\partial o_{t+1}}{\partial h_{t+1}} \cdot \sum_{k=1}^{t+1} \frac{\partial h_{t+1}}{\partial h_k} \cdot \frac{\partial h_k}{\partial w_n}$$

(esto es otra regla de la cadena mas!!)

$$\frac{\partial h_{t+1}}{\partial h_k} = \prod_{j=k}^t \frac{\partial h_{j+1}}{\partial h_j} \rightarrow \frac{\partial h_{t+1}}{\partial h_k} = \frac{\partial h_{t+1}}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_{t-1}} \cdot \dots \cdot \frac{\partial h_{k+1}}{\partial h_k}$$

algunos términos de la cadena

Ahora si unimos todo

$$\frac{\partial L_{t+1}}{\partial w_n} = \frac{\partial L_{t+1}}{\partial \hat{y}_{t+1}} \cdot \frac{\partial \hat{y}_{t+1}}{\partial o_{t+1}} \cdot \frac{\partial o_{t+1}}{\partial h_{t+1}} \cdot \sum_{k=1}^{t+1} \left[\prod_{j=k}^t \left(\frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_k}{\partial w_n} \right]$$

obd. es para 1 solo tiempo!

$$\frac{\partial L_T}{\partial W_h} = \sum_t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial h_t} \cdot \sum_{k=1}^t \left[\prod_{i=k}^t \left(\frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W_h} \right] \quad \text{V}$$

1 x rec.



9

4º Para $\frac{\partial L_T}{\partial W_x}$ se computa igual que 3º ($\frac{\partial L_T}{\partial W_h}$)



Se repite cada $h_{t-1} = 1$ a función de W_x y h_{t-1}

$$\frac{\partial L_T}{\partial W_x} = \sum_t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial h_t} \cdot \sum_{k=1}^t \left[\prod_{i=k}^t \left(\frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W_x} \right]$$

$$h_t = x_t W_x + h_{t-1} W_h + b_h$$

$$h_{t+1} = x_{t+1} W_x + h_t W_h + b_h$$

de acá cambia

5º Igual procedimiento

$$\frac{\partial L_T}{\partial b_h} = \sum_t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial h_t} \cdot \sum_{k=1}^t \left[\prod_{i=k}^t \left(\frac{\partial h_i}{\partial h_{i-1}} \right) \cdot \frac{\partial h_k}{\partial b_h} \right]$$

↓
= 1 de ecuación

5 nota I