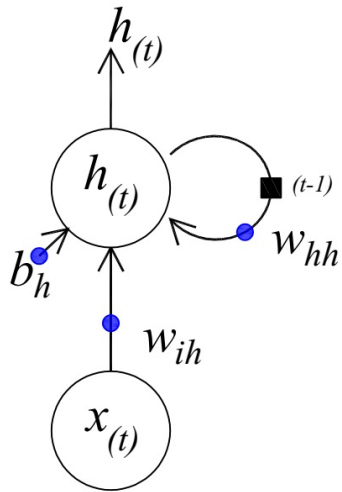


Back Propagation Through Time (BPTT)

Modelo a emplear:

$$h(t) = \tanh \left(\underbrace{x_t \cdot w_{ih} + h_{t-1} w_{hh} + b_h}_{a(t)} \right)$$



Parámetros a optimizar:

1º w_{hh} → Pesos hidden-hidden

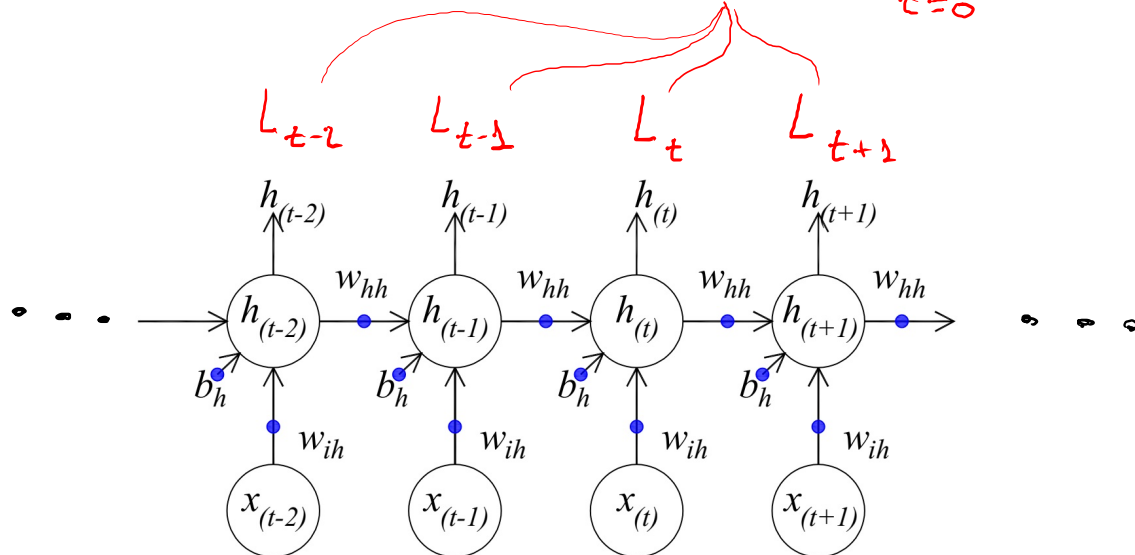
2º w_{ih} → Pesos input-hidden

3º b_h → bias hidden

¿cómo varía mi costo (L) respecto de ellos?

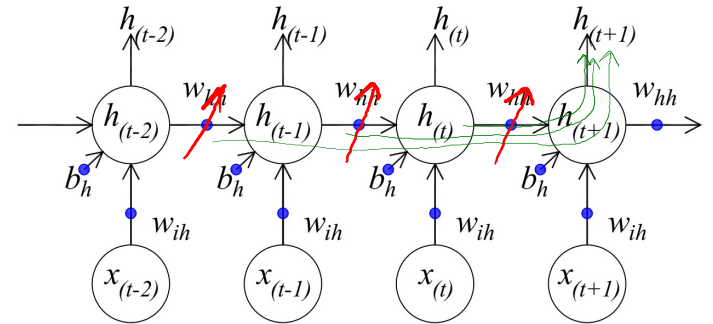
Si representamos la versión "unfolded"

costo total → $L_T = \sum_{t=0}^T L_t$



1º $w_{hh} \rightarrow$ pesos hidden-hidden

$$\frac{\partial L_T}{\partial w_{hh}} = ?$$



(A)

$$\frac{\partial L_{t+1}}{\partial w_{hh}} = \underbrace{\frac{\partial L_{t+1}}{\partial L_F} \cdot \frac{\partial L_F}{\partial h_{t+1}}}_{\text{derivada de Loss Function respecto a } h_{t+1}} \cdot \underbrace{\frac{\partial h_{t+1}}{\partial w_{hh}}}_{\text{A RESOLVER!}}$$

derivada de Loss Function
respecto a h_{t+1}

A RESOLVER!

$$h_{t+1} = \tanh \left(x_{t+1} \cdot w_{ih} + h_{t-1} w_{hh} + b_h \right)$$

h_{t+1} es función de w_{hh} y h_t

que es función de
 w_{hh} y h_{t-1}
...

$$\frac{\partial h_{t+1}}{\partial w_{hh}} = \frac{\partial h_{t+1}}{\partial h_t} \cdot \frac{\partial h_t}{\partial w_{hh}}$$

se propone escribirlo como

$$\frac{\partial h_{t+1}}{\partial w_{hh}} = \sum_{k=0}^{t+1} \frac{\partial h_{t+1}}{\partial h_k} \cdot \frac{\partial h_k}{\partial w_{hh}}$$

ejemplo

$$\frac{\partial h_2}{\partial w_{hh}} = \sum_{k=0}^2 \frac{\partial h_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial w_{hh}} =$$

$$= \underbrace{\frac{\partial h_2}{\partial h_0} \frac{\partial h_0}{\partial w_{hh}}}_{h_2 \text{ varia porque } h_0 \text{ varia por que } w_{hh} \text{ varia}} + \underbrace{\frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_{hh}}}_{h_2 \text{ varia porque } h_1 \text{ varia por que } w_{hh} \text{ varia}} + \underbrace{\frac{\partial h_2}{\partial h_2} \frac{\partial h_2}{\partial w_{hh}}}_{h_2 \text{ varia porque } w_{hh} \text{ varia}}$$

h_2 varia porque h_0 varia por que w_{hh} varia
 h_2 varia porque h_1 varia por que w_{hh} varia
 h_2 varia porque w_{hh} varia

Retomamos (A)

$$\frac{\partial L_{t+1}}{\partial w_{hh}} = \frac{\partial L_{t+1}}{\partial LF} \cdot \frac{\partial LF}{\partial h_{t+1}} \cdot \sum_{k=0}^{t+1} \underbrace{\frac{\partial h_{t+1}}{\partial h_k}}_{\text{ES OTRA REGLA DE LA CADENA !}} \frac{\partial h_k}{\partial w_{hh}}$$

ES OTRA
REGLA DE
LA CADENA !

$$\frac{\partial h_{t+1}}{\partial h_k} = \prod_{j=k}^t \frac{\partial h_{j+1}}{\partial h_j} \rightarrow \underbrace{\frac{\partial h_{k+1}}{\partial h_k}}_{\text{el 1º término}} \cdot \dots \cdot \underbrace{\frac{\partial h_{t+1}}{\partial h_t}}_{\text{el último término}}$$

Unimos todo...

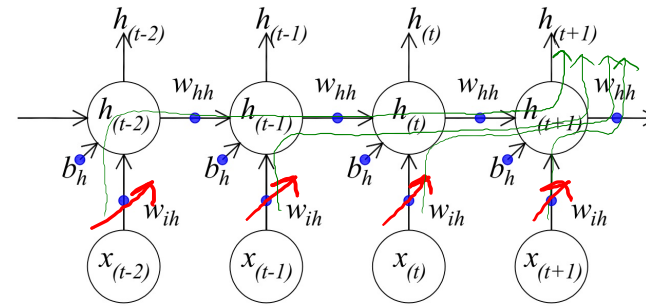
$$\frac{\partial L_{t+1}}{\partial w_{hh}} = \frac{\partial L}{\partial LF} \cdot \frac{\partial LF}{\partial h_{t+1}} \cdot \sum_{k=0}^{t+1} \left[\prod_{j=k}^{t+1} \left(\frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_k}{\partial w_{hh}} \right]$$

Solo para $t+1$!!

$$\frac{\partial L_T}{\partial w_{hh}} = \sum_{t=0}^T \frac{\partial L_t}{\partial LF} \cdot \frac{\partial LF}{\partial h_t} \sum_{k=0}^t \left[\prod_{j=k}^t \left(\frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_k}{\partial w_{hh}} \right]$$

2º w_{ih} → Pesos input-hidden

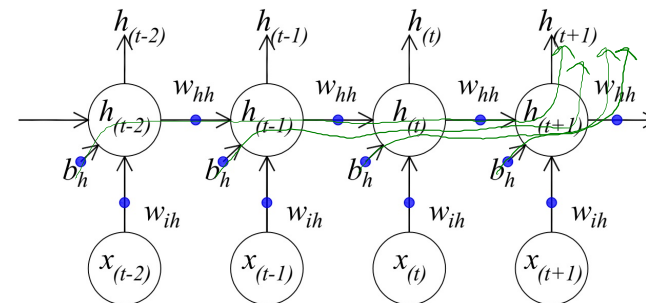
Procedimiento semejante ya que h_t es función de w_{ih} y h_{t-1}



$$\frac{\partial L_T}{\partial w_{ih}} = \sum_{t=0}^T \frac{\partial L_t}{\partial L_F} \cdot \frac{\partial L_F}{\partial h_t} \sum_{k=0}^t \left[\prod_{j=k}^t \left(\frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_k}{\partial w_{ih}} \right]$$

acá cambió

3º b_h → bias hidden
idem



$$\frac{\partial L_T}{\partial b_h} = \sum_{t=0}^T \frac{\partial L_t}{\partial L_F} \cdot \frac{\partial L_F}{\partial h_t} \sum_{k=0}^t \left[\prod_{j=k}^t \left(\frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_k}{\partial b_h} \right]$$

acá cambió