

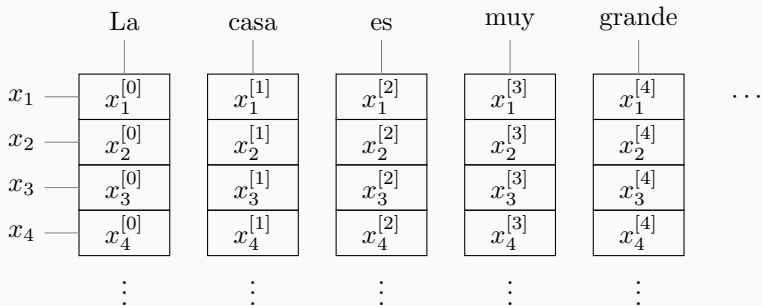
Aprendizaje profundo

REDES RECURRENTE

Gibran Fuentes-Pineda

Octubre 2021

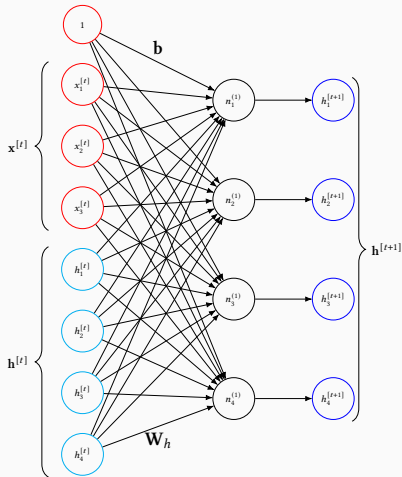
Motivación: secuencias de palabras



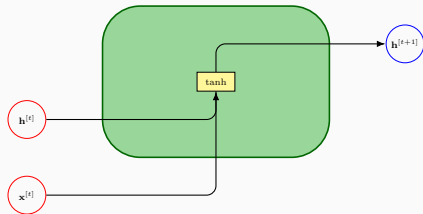
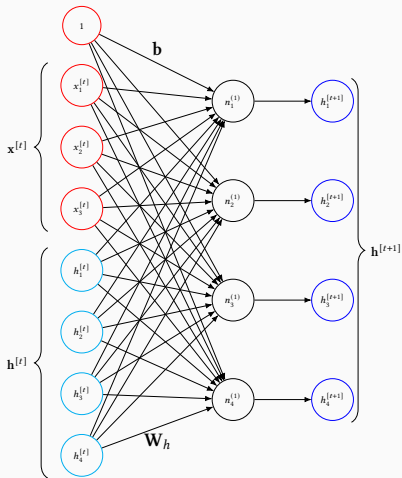
Unidad recurrente básica

- Capas con retro-alimentación en sus conexiones
 1. Entradas en tiempo t ($\mathbf{x}^{[t]}$)
 2. Estado en tiempo t ($\mathbf{h}^{[t]}$)

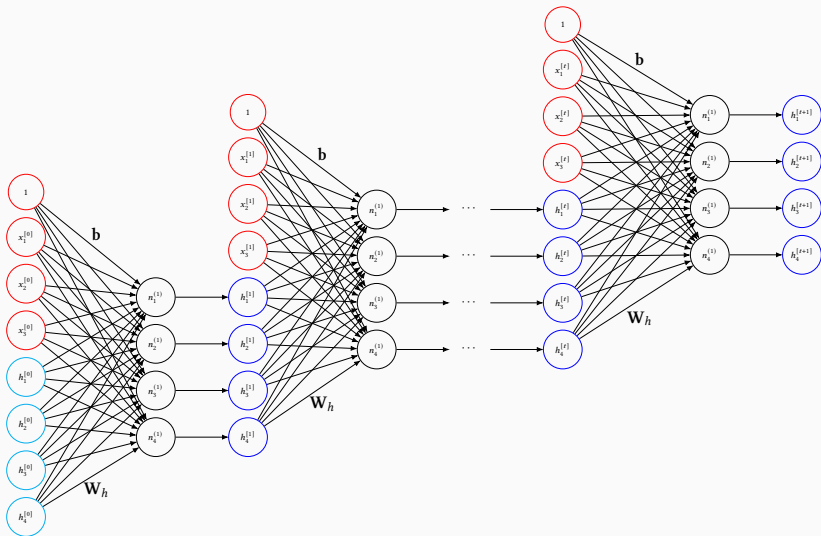
$$\begin{aligned} \mathbf{h}^{[t+1]} &= \phi \left(\mathbf{w}_h \cdot \underbrace{\begin{bmatrix} \mathbf{h}^{[t]}, \mathbf{x}^{[t]} \end{bmatrix}}_{\text{Concatenación}} + \mathbf{b}_h \right) \\ &= \phi \left(\mathbf{w}_{hh} \cdot \mathbf{h}^{[t]} + \mathbf{w}_{hx} \cdot \mathbf{x}^{[t]} + \mathbf{b}_h \right) \end{aligned}$$



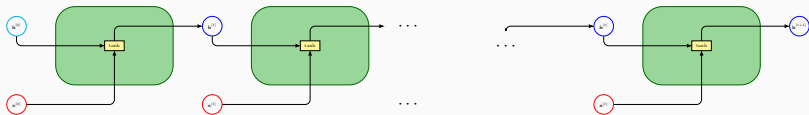
Unidad recurrente básica: diagrama de celda



Unidad recurrente básica: despliegue

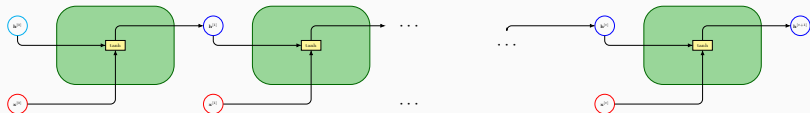


Unidad recurrente básica: despliegue de celdas



Modelando dependencias a corto plazo

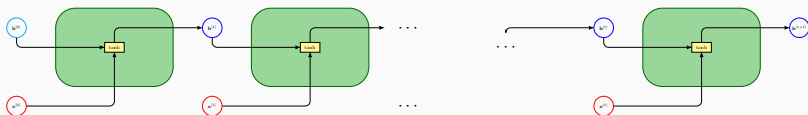
- En teoría una red recurrente básica puede modelar dependencias a corto y largo plazo
 - Siegelmann y Sontag mostraron que las redes recurrentes son Turing completas¹



¹Siegelmann and Sontag. On The Computational Power Of Neural Nets, 1995.

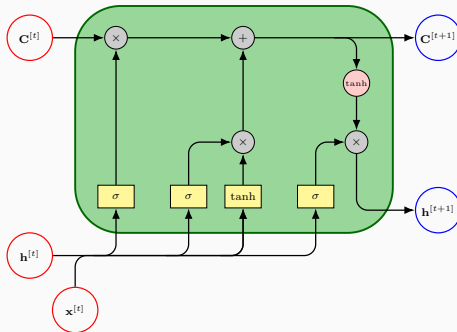
El problema de la memoria a largo plazo

- En práctica es muy difícil entrenarlas para tareas con dependencias a largo plazo debido al problema del desvanecimiento/explosión de los gradientes



Long-short term memory (LSTM)²

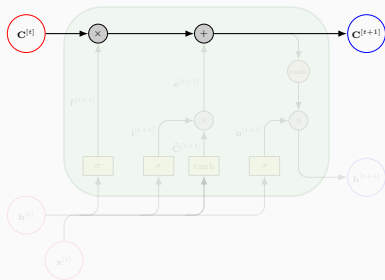
- Agregan elementos internos a la celda básica que permiten capturar dependencias a corto y largo plazo



²Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*. 9 (8): 1735–1780, 1997.

Long-short term memory (LSTM): salida de la capa anterior

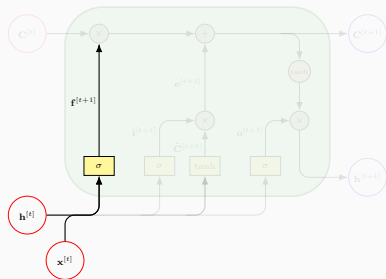
- Agrega o elimina información del estado $\mathbf{C}^{[t]}$ a partir de la entrada actual $\mathbf{x}^{[t+1]}$ y la salida anterior $\mathbf{h}^{[t]}$



Long-short term memory (LSTM): compuerta de olvido

- Determina qué olvidar del estado $C^{[t]}$ y en qué proporción a partir de la entrada actual $x^{[t+1]}$ y la salida anterior $h^{[t]}$

$$f^{[t+1]} = \sigma \left(W_f \cdot \begin{bmatrix} h^{[t]} \\ x^{[t+1]} \end{bmatrix} + b_f \right)$$



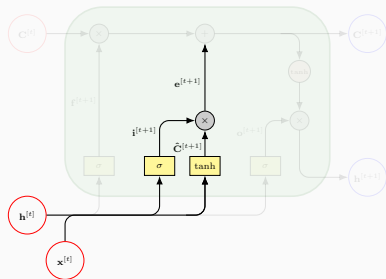
Long-short term memory (LSTM): computerta de entrada

- Determina qué agregar al estado $\mathbf{C}^{[t]}$ y en qué proporción a partir de la entrada actual $\mathbf{x}^{[t+1]}$ y el estado oculto anterior $\mathbf{h}^{[t]}$

$$\mathbf{i}^{[t+1]} = \sigma \left(\mathbf{W}_i \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_i \right)$$

$$\hat{\mathbf{C}}^{[t+1]} = \tanh \left(\mathbf{W}_c \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_c \right)$$

$$\mathbf{e}^{[t+1]} = \mathbf{i}^{[t+1]} \odot \hat{\mathbf{C}}^{[t+1]}$$

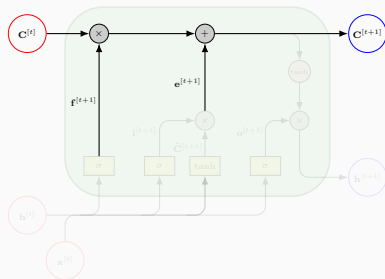


Long-short term memory (LSTM): nuevo estado

- El nuevo estado $\mathbf{C}^{[t+1]}$ se obtiene como una combinación del estado $\mathbf{C}^{[t]}$, la salida $\mathbf{f}^{(t)}$ de la compuerta de olvido y la salida $\mathbf{e}^{[t+1]}$ de la compuerta de entrada

$$\mathbf{C}^{[t+1]} = \mathbf{f}^{[t+1]} \odot \mathbf{C}^{[t]} + \mathbf{e}^{[t+1]}$$

donde \odot denota el producto de Hadamard

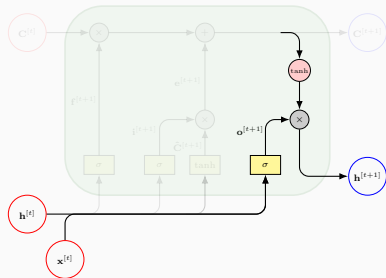


Long-short term memory (LSTM): computerta de salida

- El siguiente estado oculto $\mathbf{h}^{[t+1]}$ se obtiene como una combinación de la entrada actual $\mathbf{x}^{[t+1]}$, el estado oculto anterior $\mathbf{h}^{[t]}$ y el nuevo estado $\mathbf{c}^{[t+1]}$

$$\mathbf{o}^{[t+1]} = \sigma \left(\mathbf{W}_o \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_o \right)$$

$$\mathbf{h}^{[t+1]} = \mathbf{o}^{[t+1]} \odot \tanh \left(\mathbf{c}^{[t+1]} \right)$$



- Combina compuertas de olvido y entrada en una sola

$$\mathbf{z}^{[t+1]} = \sigma \left(\mathbf{W}_z \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_z \right)$$

$$\mathbf{r}^{[t+1]} = \sigma \left(\mathbf{W}_r \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_r \right)$$

$$\tilde{\mathbf{h}}^{[t+1]} = \tanh \left(\mathbf{W}_h \cdot \left[\mathbf{r}^{[t+1]} \odot \mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_h \right)$$

$$\mathbf{h}^{[t+1]} = \left(1 - \mathbf{z}^{[t+1]} \right) \odot \mathbf{h}^{[t]} + \mathbf{z}^{[t+1]} \odot \tilde{\mathbf{h}}^{[t+1]}$$

³K. Cho et al. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *arXiv:1406.1078*, 2014.

- Contienen celdas recurrentes en conjunto con otras capas
- La salida de una celda alimenta otras capas u otras celdas
- Por ejemplo, para predecir el siguiente símbolo en un texto con una celda recurrente básica, a la salida podemos agregar una capa densa con función de activación *softmax*

$$\hat{y}^{[t+1]} = \text{softmax} \left(\mathbf{w}_y \cdot \mathbf{h}^{[t+1]} + \mathbf{b}_y \right)$$

Arquitecturas de redes recurrentes: ejemplo

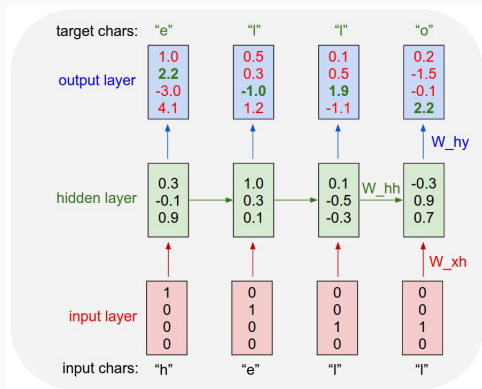


Imagen tomada de Karpathy 2015 (<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>)

Arquitecturas de redes recurrentes: tareas de uno a uno

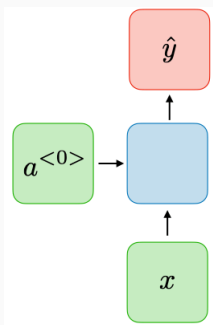


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de uno a muchos

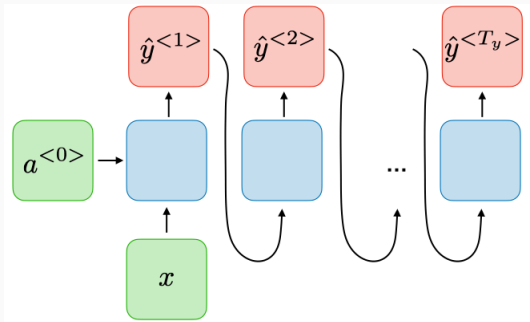


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de muchos a uno

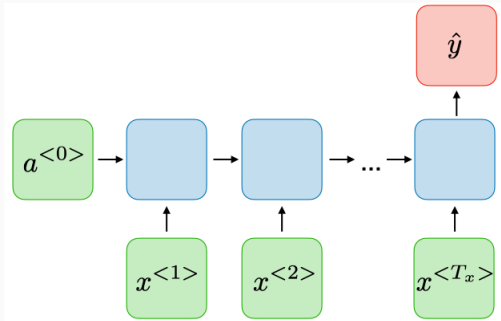
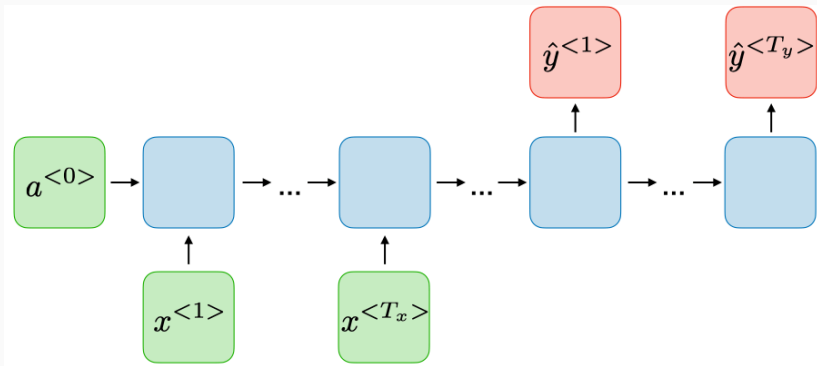


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de muchos a muchos



Imagen

tomada de Amidi, Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: LSTM/GRU bidireccional

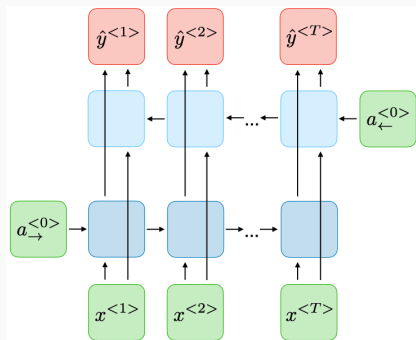


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: celdas apiladas

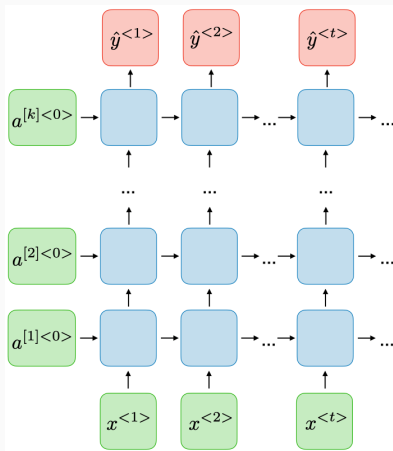


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

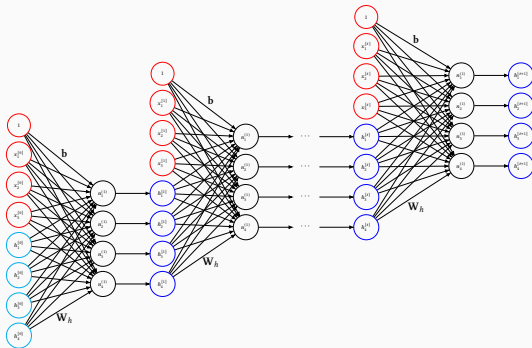
Retropropagación en el tiempo

- Pérdida en el tiempo

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^T \mathcal{L}(\hat{y}^{[t]}, y^{[t]})$$

- Retropropagación

$$\frac{\partial \mathcal{L}^{[T]}}{\partial \theta} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{[t]}}{\partial \theta}$$



Retropropagación en el tiempo para una celda básica (1)

- Para la matriz de pesos \mathbf{W}_y y un tiempo T

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_y} = \sum_{t=1}^T \left[\frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[t]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[t]}}{\partial \mathbf{W}_y} \right] = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[t]}} \mathbf{h}^{[t]\top}$$
$$\frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[t]}} = \frac{1}{T} \cdot \frac{\partial \mathcal{L}(\hat{\mathbf{y}}^{[t]}, \mathbf{y}^{[t]})}{\partial \hat{\mathbf{y}}^{[t]}}$$

- Para el tiempo T

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[T]}} = \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[T]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[T]}}{\partial \mathbf{h}^{[T]}} = \mathbf{W}_y^\top \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[T]}}$$

Retropropagación en el tiempo para una celda básica (2)

- Para los tiempos $T - 1, \dots, 1$, la función de pérdida se ve afectada por $\mathbf{h}^{[t]}$ a través de $\mathbf{h}^{[t+1]}$ y $\hat{\mathbf{y}}^{[t]}$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[t]}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[t+1]}} \cdot \frac{\partial \mathbf{h}^{[t+1]}}{\partial \mathbf{h}^{[t]}} + \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[t]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[t]}}{\partial \mathbf{h}^{[t]}} \\ &= \mathbf{W}_{hh}^\top \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[t+1]}} + \mathbf{W}_y^\top \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}^{[t]}}\end{aligned}$$

- Para la matriz de pesos \mathbf{W}_{hh} y un tiempo T

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hh}} &= \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}} \cdot \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}^{[t]}} \cdot \left[\sum_{t=1}^T \frac{\partial \mathbf{h}^{[T]}}{\partial \mathbf{h}^{[t]}} \cdot \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{W}_{hh}} \right] \\ \frac{\partial \mathbf{h}^{[T]}}{\partial \mathbf{h}^{[t]}} &= \prod_{i=t+1}^T \frac{\partial \mathbf{h}^{[i]}}{\partial \mathbf{h}^{[i-1]}}\end{aligned}$$

Retropropagación en el tiempo para una celda básica (3)

- Para la matriz de pesos \mathbf{W}_{hx} y un tiempo T

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hx}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[t]}} \cdot \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{W}_{hx}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\mathbf{h}^{[t]}} \cdot \mathbf{x}^{[t]\top}$$

- Para la matriz de pesos \mathbf{W}_{hh} y un tiempo T

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hh}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{[t]}} \cdot \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{W}_{hh}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\mathbf{h}^{[t]}} \cdot \mathbf{h}^{[t-1]\top}$$