

Aprendizaje profundo

REDES RECURRENTES

Gibran Fuentes-Pineda

Noviembre 2020

Tareas uno a uno

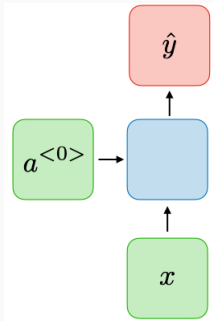


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Tareas uno a muchos

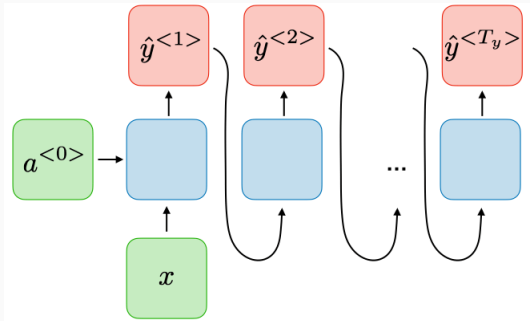


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Tareas muchos a uno

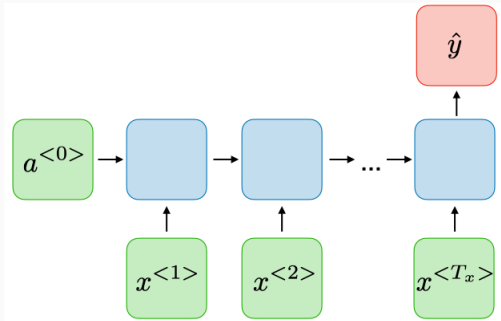


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Tareas muchos a muchos

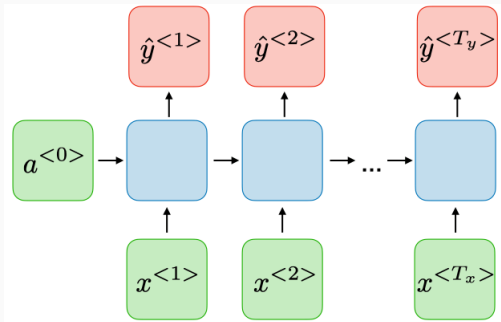


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Tareas muchos a muchos con tiempos distintos

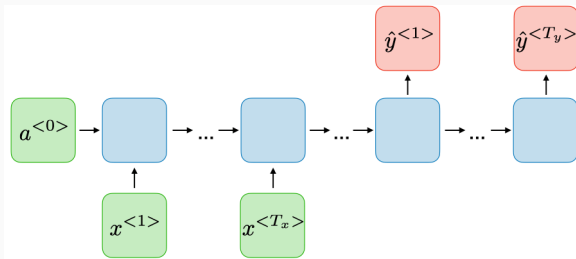


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

- Redes con retro-alimentación en sus conexiones

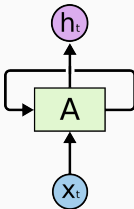
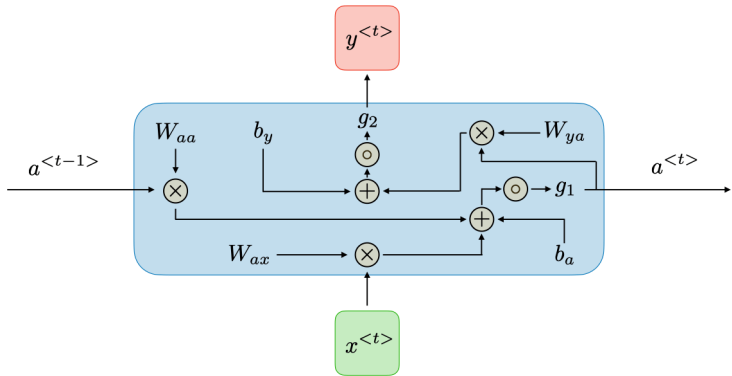


Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

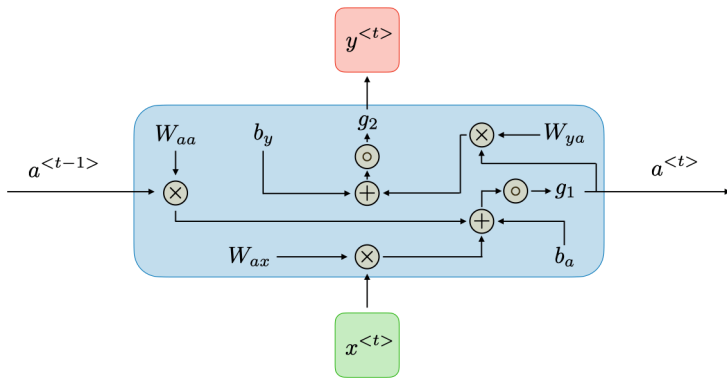
Elementos básicos

1. Entradas en cada instante de tiempo t ($\mathbf{x}^{[t+1]}$)
2. Estados en cada instante de tiempo t ($\mathbf{h}^{[t+1]}$)
3. Salidas en cada instante de tiempo t ($\mathbf{y}^{[t+1]}$)



Tipos de retroalimentación

1. Salida anterior ($y^{[t]}$)
2. Estado anterior ($h^{[t]}$)
3. Salida y estado anterior ($h^{[t]}$ y $y^{[t]}$)



Celda recurrente básica

- Capa que procesa el estado anterior y la entrada actual para generar un nuevo estado y la salida

$$\mathbf{h}^{[t+1]} = \phi(\mathbf{W}_{hh} \cdot \mathbf{h}^{[t]} + \mathbf{W}_{hx} \cdot \mathbf{x}^{[t+1]} + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}^{[t+1]} = \phi(\mathbf{W}_{yh} \cdot \mathbf{h}^{[t+1]} + \mathbf{b}_y)$$

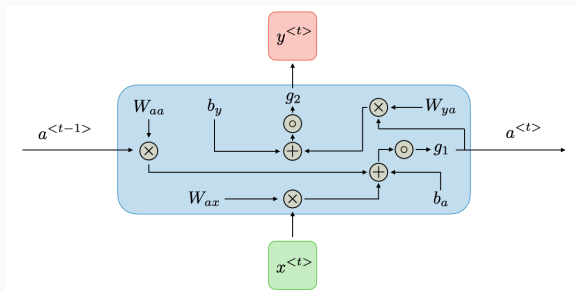
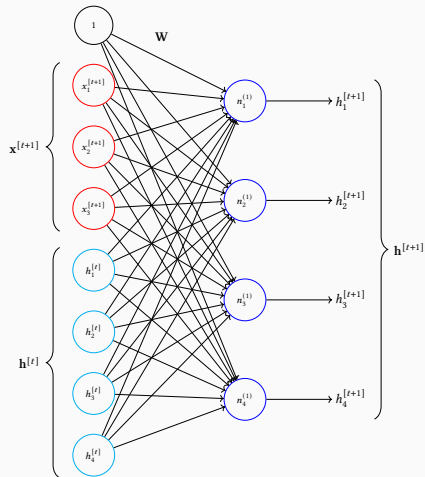


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Celda recurrente básica: otra perspectiva

$$\mathbf{h}^{[t+1]} = \sigma \left(\mathbf{W}_h \cdot \underbrace{\begin{bmatrix} \mathbf{h}^{[t]} \\ \mathbf{x}^{[t+1]} \end{bmatrix}}_{\text{Concatenación}} + \mathbf{b}_h \right)$$



Despliegue de celdas

- Una celda recurrente para una secuencia de T valores, se puede desplegar en T capas con parámetros idénticos

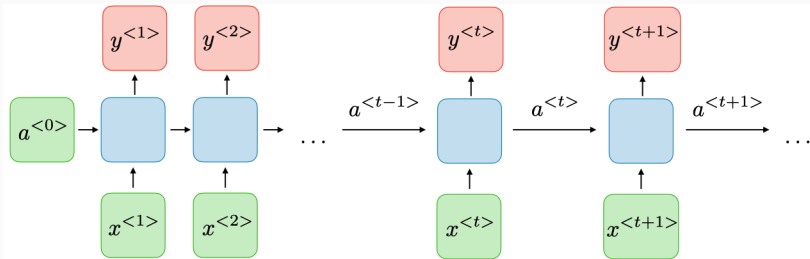


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Red neuronal con celdas recurrentes

- Típicamente contiene celdas recurrentes en conjunto con otras capas
- La salida de una celda puede alimentar otras capas u otras celdas
- Un clasificador simple

$$\mathbf{h}^{[t+1]} = \sigma \left(\mathbf{w}_h \cdot \left[\mathbf{h}^{[t]}, \mathbf{x}^{[t+1]} \right] + \mathbf{b}_h \right)$$
$$\hat{\mathbf{y}}^{[t+1]} = \textit{softmax} \left(\mathbf{w}_{yh} \cdot \mathbf{h}^{[t+1]} + \mathbf{b}_y \right)$$

Ejemplo: modelo de lenguaje a nivel símbolo

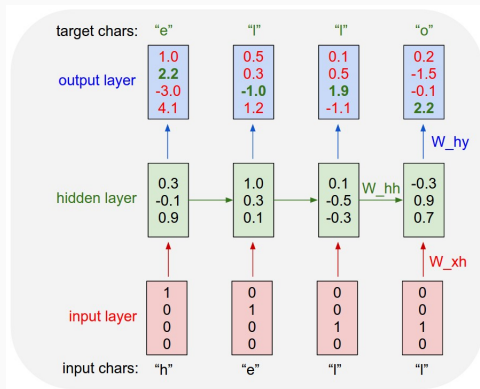


Imagen tomada de Karpathy 2015 (<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>)

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

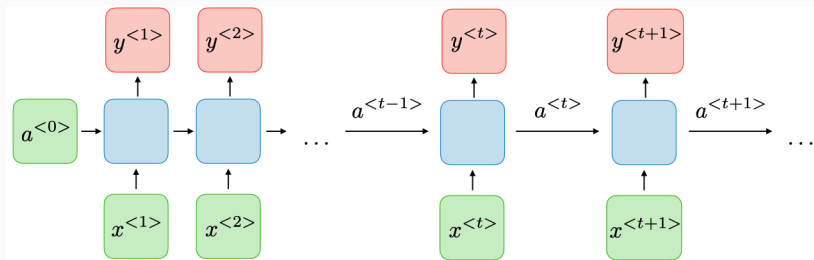


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

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

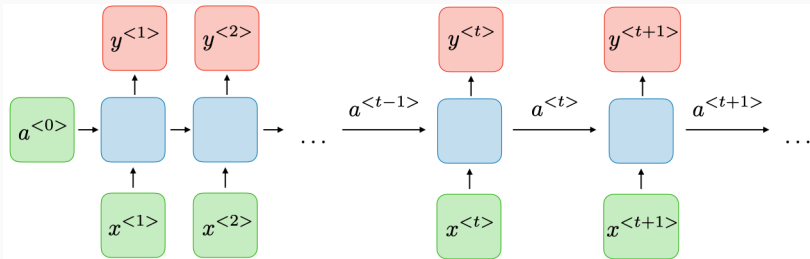


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Memorias a corto y largo plazo

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

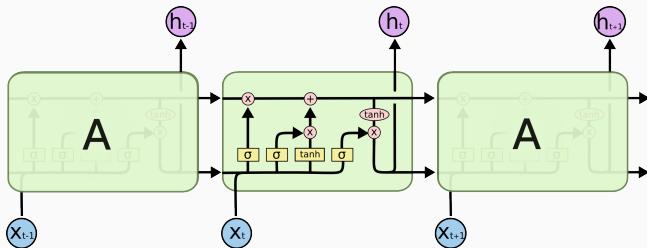


Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

LSTM: salida de la capa anterior

- Agrega o elimina elementos del estado anterior de la celda $C^{[t]}$ basado en transformación de la entrada actual $x^{[t+1]}$ y el estado oculto anterior $h^{[t]}$

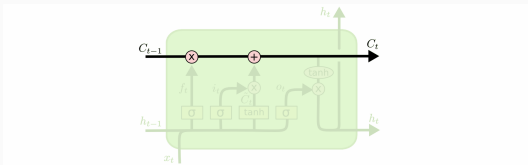


Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

LSTM: compuerta de olvido

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

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

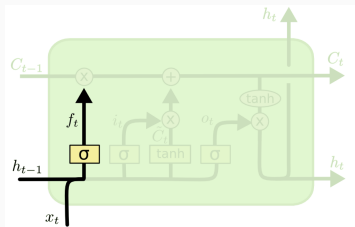


Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

LSTM: computerta de entrada

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

$$i^{[t+1]} = \sigma \left(W_i \cdot \left[h^{[t]}, x^{[t+1]} \right] + b_i \right)$$

$$\hat{C}^{[t+1]} = \tanh \left(W_C \cdot \left[h^{[t]}, x^{[t+1]} \right] + b_C \right)$$

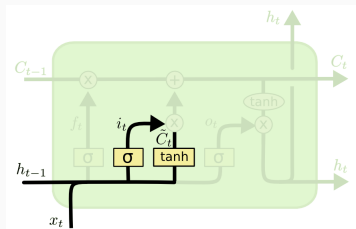


Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

LSTM: nuevo estado

- El nuevo estado de la celda se obtiene como una combinación de la salida de la compuerta de olvido $f^{(t)}$ y las salidas $i^{[t+1]}$ y $\tilde{C}^{[t+1]}$ de la compuerta de entrada

$$C^{[t+1]} = f^{[t+1]} \odot C^{[t]} + i^{[t+1]} \odot \tilde{C}^{[t+1]}$$

donde \odot denota el producto de Hadamard

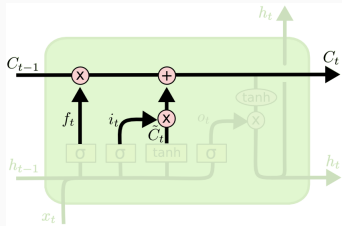


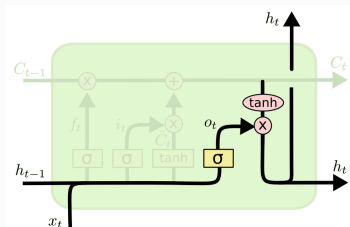
Imagen tomada de Colah 2015 (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

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 de la celda $\mathbf{C}^{[t+1]}$

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

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



Gated Recurrent Unit

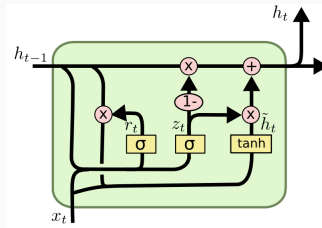
- Combina compuertas de olvido y entrada en una sola

$$z^{[t+1]} = \sigma \left(W_z \cdot \left[h^{[t]}, x^{[t+1]} \right] + b_z \right)$$

$$r^{[t+1]} = \sigma \left(W_r \cdot \left[h^{[t]}, x^{[t+1]} \right] + b_r \right)$$

$$\tilde{h}^{[t+1]} = \tanh \left(W_h \cdot \left[r^{[t+1]} \odot h^{[t]}, x^{[t+1]} \right] + b_h \right)$$

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



Retropropagación en el tiempo

- Pérdida en el tiempo

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

- Retropropagación

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

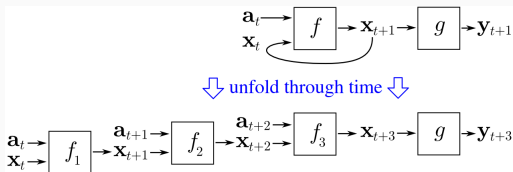


Imagen tomada de Wikipedia (https://en.wikipedia.org/wiki/Backpropagation_through_time)

Redes recurrentes apiladas

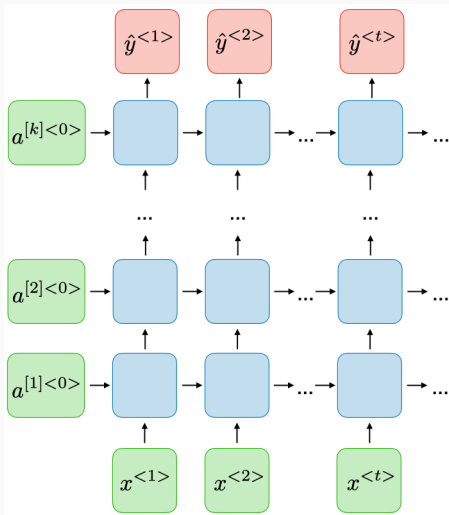


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

RNR Bidireccional

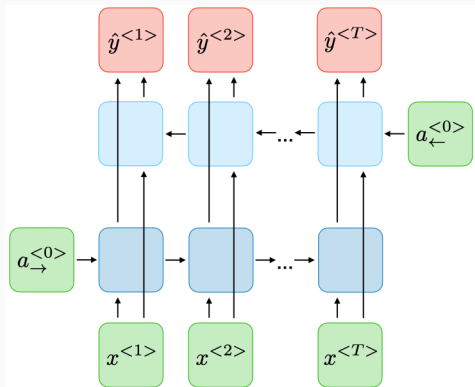


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Modelos secuencia a secuencia

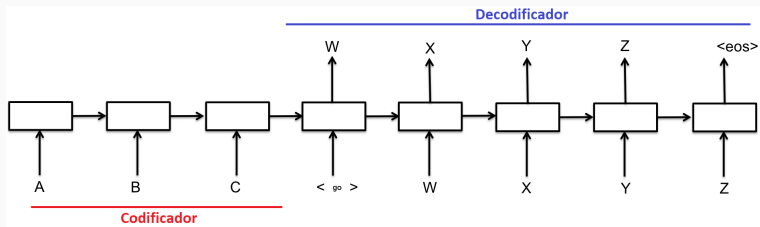


Imagen derivada de <https://www.tensorflow.org/tutorials/seq2seq>