

Aprendizaje profundo

REDES RECURRENTES

Gibran Fuentes-Pineda

Octubre 2023

- 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

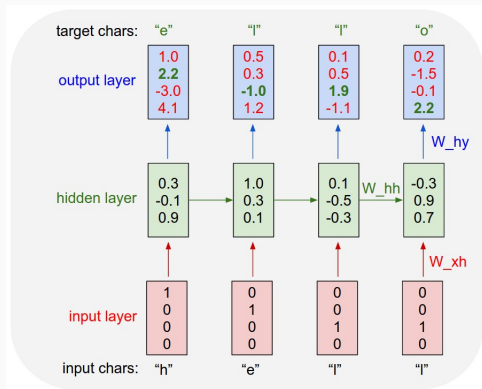


Imagen tomada de Karpathy 2015 (<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>)

Arquitecturas de redes recurrentes: tareas de uno a uno

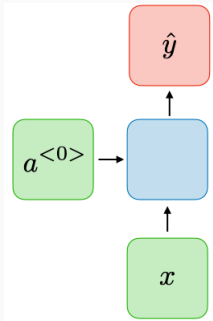


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de uno a muchos

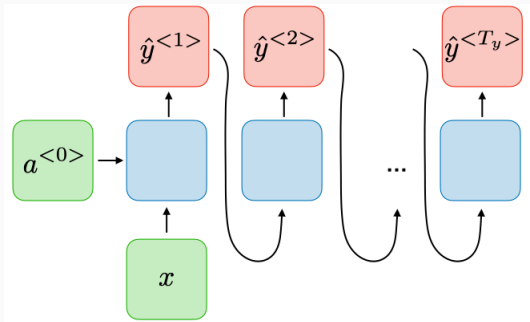


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de muchos a uno

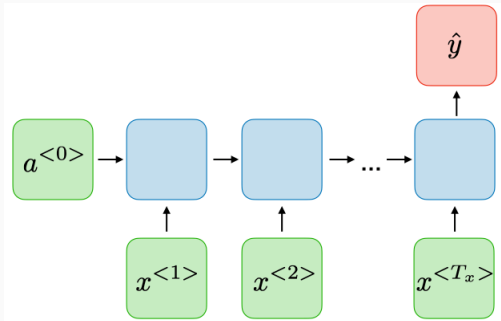
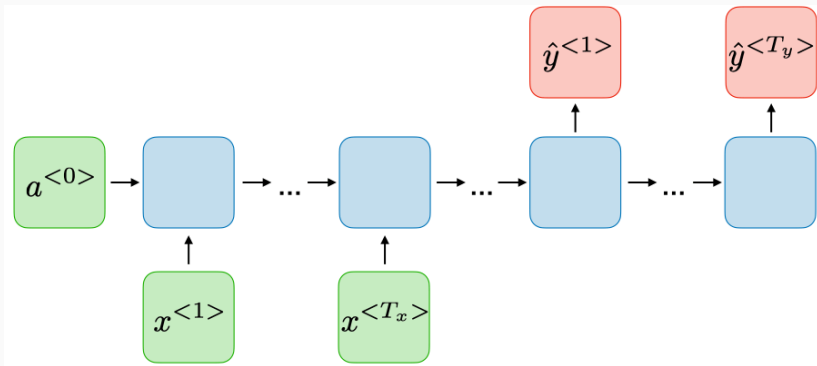


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: tareas de muchos a muchos



Imagen

tomada de Amidi, Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: LSTM/GRU bidireccional

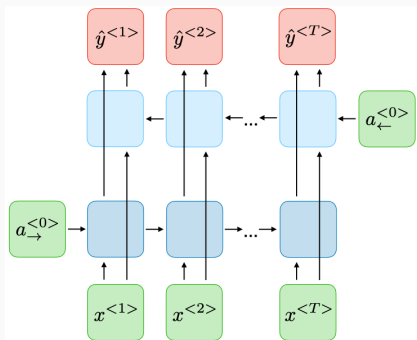


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

Arquitecturas de redes recurrentes: celdas apiladas

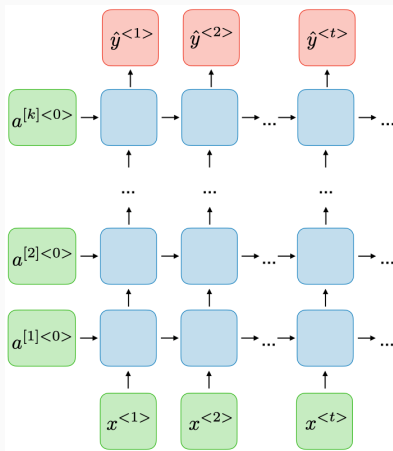


Imagen tomada de Amidi. Recurrent Neural Networks cheatsheet

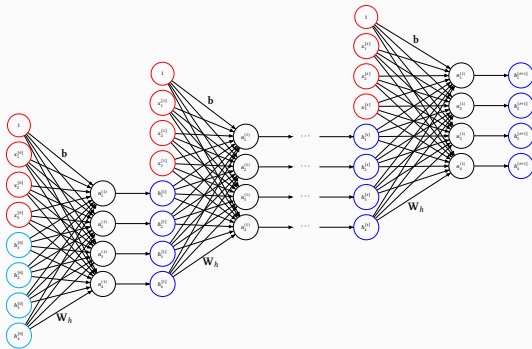
Retropropagación en el tiempo

- Pérdida en el tiempo

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

- Retropropagación

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



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

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

Retropropagación en el tiempo para una celda básica (3)

- La pérdida depende de \mathbf{W}_{hx} y \mathbf{W}_{hh} por $\mathbf{h}^{[1]}, \mathbf{h}^{[2]}, \dots, \mathbf{h}^{[T]}$
- Para la matriz de pesos \mathbf{W}_{hx}

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

- Para la matriz de pesos \mathbf{W}_{hh}

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