

# **MAC5921 – Deep Learning**

Aula 11 – 26/09/2023

## **RNN**

Nina S. T. Hirata

Usadas para dados cujos elementos são organizados em uma estrutura sequencial

- série temporal
- sequência de genes
- áudio
- vídeo
- texto

## Redes neurais convolucionais

- explora o fato de um mesmo padrão poder ocorrer em qualquer parte de uma imagem
- convolução considera uma vizinhança restrita; múltiplas camadas para “ampliar” a vizinhança abarcada — relação espacial
- o compartilhamento de pesos (entre diferentes posições na imagem) permite a aplicação em imagens de variados tamanhos

## Redes neurais recorrentes

- explora o fato de um mesmo padrão poder ocorrer em qualquer parte de uma sequência
- recorrência considera todos os elementos prévios na sequência; não considera elementos posteriores na sequência — relação sequencial
- o compartilhamento de pesos permite a aplicação em sequências de variados tamanhos

## Entrada:

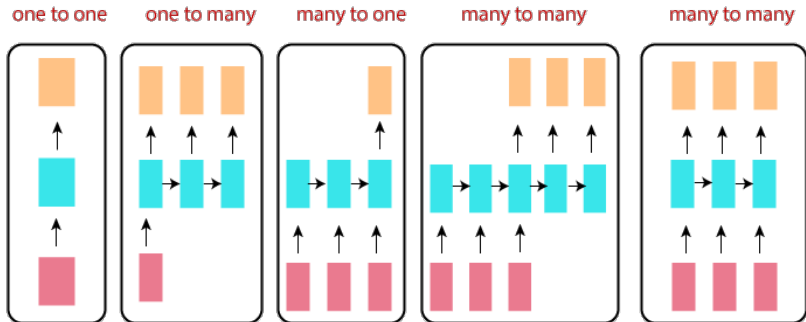
Sequência de elementos; elementos são processados um a um

$$\mathbf{x} = x_1, x_2, \dots, x_t, \dots, x_T$$

## Saída pode ser:

- scores para rótulos de classe (problema de classificação)
- sequência de elementos (por exemplo, tradução)
- scores ou valores para próximo elemento (predição de próximo elemento)
- ...

## Diferentes formatos de saída e formas de mapear entrada para saída



<https://www.javatpoint.com/tensorflow-types-of-rnn>

$t$ : índice na sequência

$x_t$ : elemento na posição  $t$  da sequência de entrada  $x$

Um novo conceito:

**Estado:**  $h_t = f(x_t, h_{t-1})$

Exemplo:  $h_t = \tanh(Ux_t + Wh_{t-1})$

Saída:  $o_t = Vh_t, \quad \hat{y}_t = \text{softmax}(Vh_t)$

$t$ : índice na sequência

$x_t$ : elemento na posição  $t$  da sequência de entrada  $x$

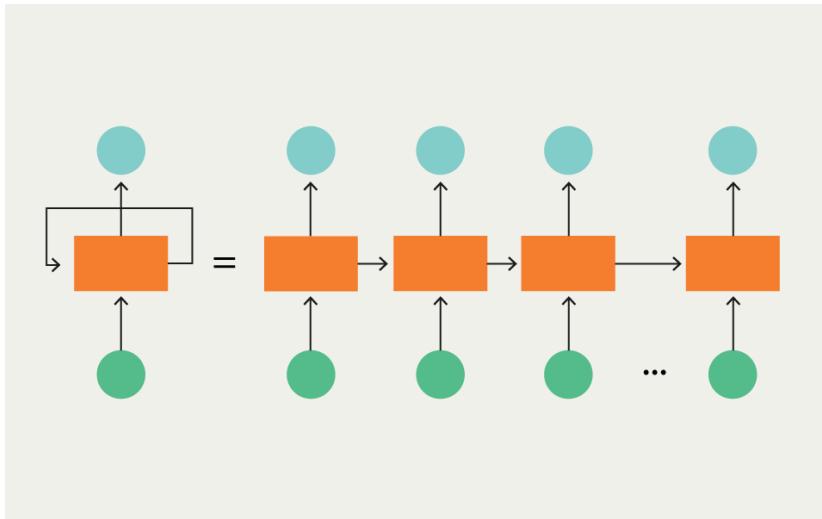
Um novo conceito:

**Estado:**  $h_t = f(x_t, h_{t-1})$

Exemplo:  $h_t = \tanh(\mathbf{U}x_t + \mathbf{W}h_{t-1})$

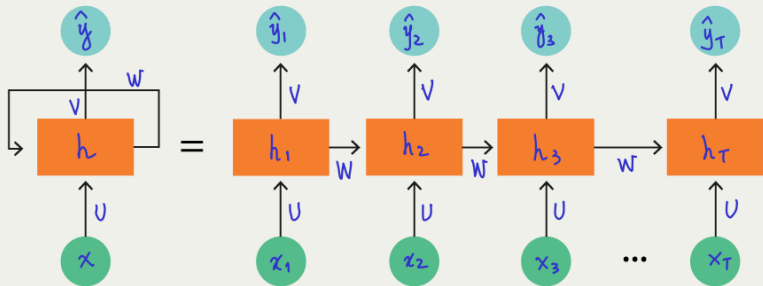
Saída:  $o_t = \mathbf{V}h_t, \quad \hat{y}_t = \text{softmax}(\mathbf{V}h_t)$

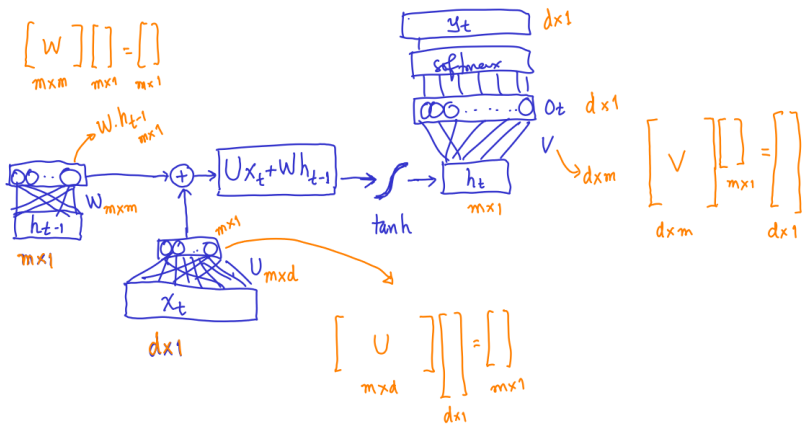
$\mathbf{U}$ ,  $\mathbf{W}$  e  $\mathbf{V}$  são matrizes de peso, compartilhados para todos os  $t$



Fonte: <https://www.telusinternational.com/articles/difference-between-cnn-and-rnn>

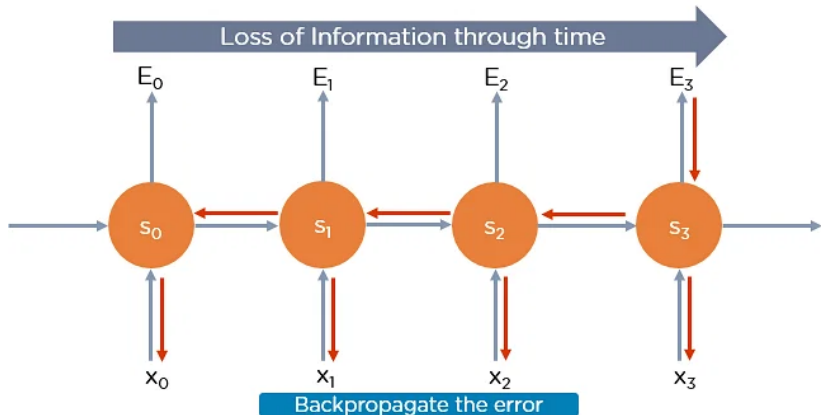




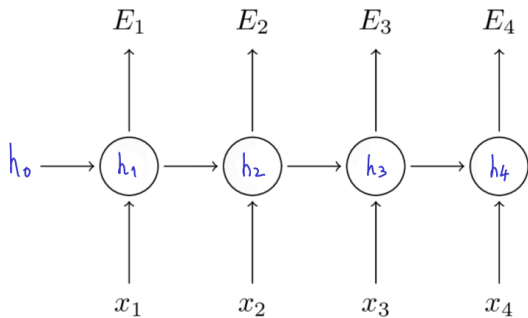


Aqui, podemos pensar que  $d$  é o tamanho do vocabulário (*one-hot-encoding*)

## Forward e Backward pass em RNNs



Fonte: <https://www.simplilearn.com/tutorials/deep-learning-tutorial/rnn>



**Erro de predição na posição  $t$**  (lembrar que  $y_t$  é one-hot-encoded target):

$$E_t(y_t, \hat{y}_t) = - \sum_c y_t \log(\hat{y}_t)$$

**Total loss**

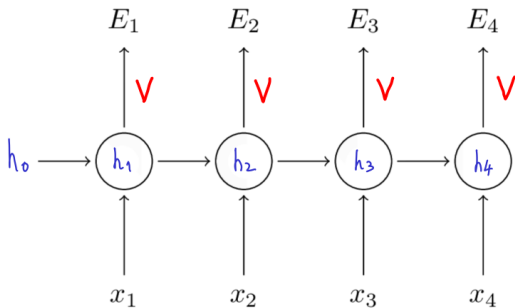
$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

**Derivadas parciais** que precisamos calcular

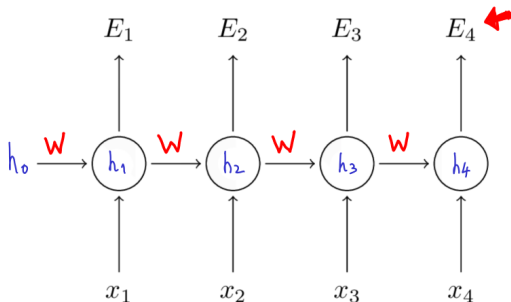
$$\frac{\partial E}{\partial u} = \sum_t \frac{\partial E_t}{\partial u}$$

$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W}$$

$$\frac{\partial E}{\partial V} = \sum_t \frac{\partial E_t}{\partial V}$$



$$\frac{\partial E}{\partial V} = \sum_t \frac{\partial E_t}{\partial V}$$



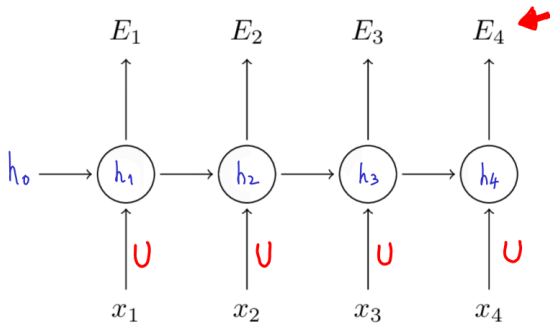
$$\frac{\partial E_t}{\partial W} = \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W} + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \dots + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_1} \frac{\partial h_1}{\partial W}$$

$$\frac{\partial E_t}{\partial W} = \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W} + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \dots + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_1} \frac{\partial h_1}{\partial W}$$

$$\begin{aligned} \frac{\partial E_t}{\partial W} &= \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W} + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \dots + \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial h_1} \frac{\partial h_1}{\partial W} \\ &= \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \left[ \frac{\partial h_t}{\partial h_t} \frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \dots + \frac{\partial h_t}{\partial h_1} \frac{\partial h_1}{\partial W} \right] \\ &= \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \sum_{i=1}^t \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W} \end{aligned}$$

Produto de gradientes pode levar a vanishing/exploding gradient





$$\frac{\partial E}{\partial u} = \sum_t \frac{\partial E_t}{\partial u}$$

Muito similar ao caso dos pesos  $W$

Exemplo de código (Python only):

<https://gist.github.com/karpathy/d4dee566867f8291f086>

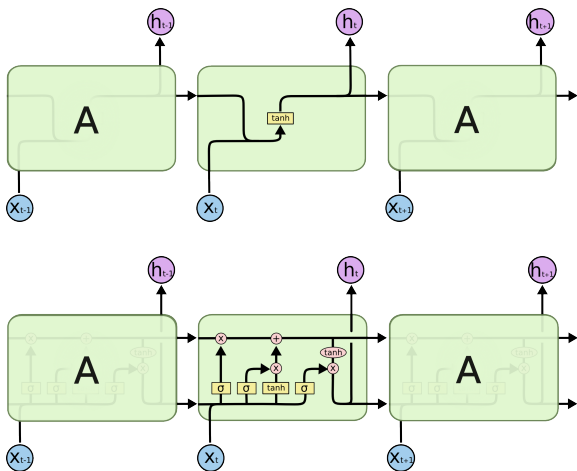
Comentário detalhado sobre o código:

<https://mkffl.github.io/2019/07/08/minimalist-RNN.html>

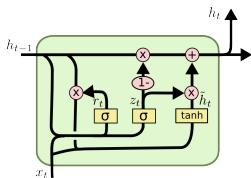
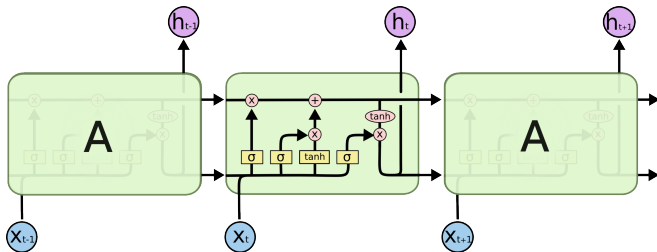
A estrutura da rede pode ser bem mais complexa

Por exemplo, múltiplas camadas ocultas

## Padrão e LSTM



# LSTM e GRU



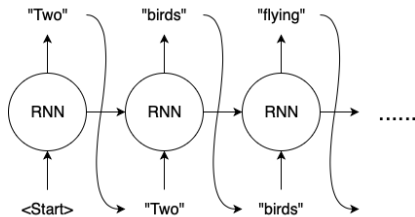
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

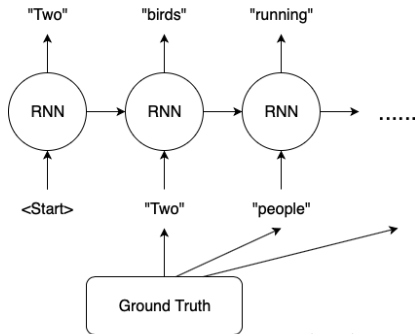
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Teacher forcing no treinamento

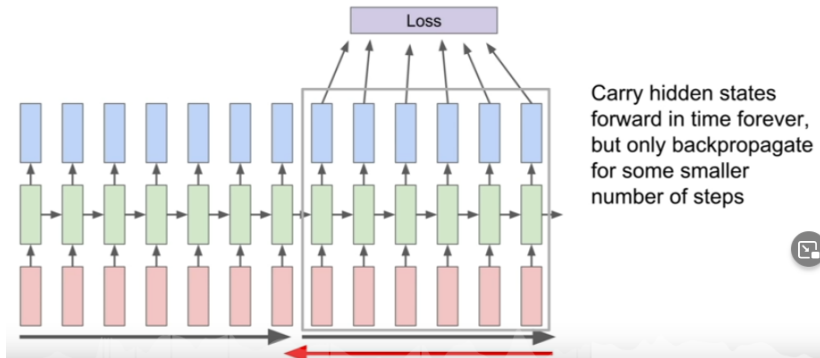


Without Teacher Forcing

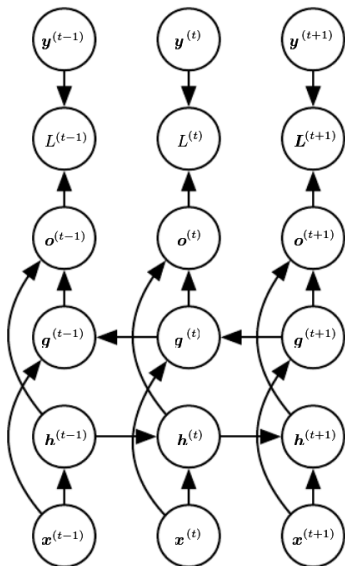


## Truncated BPTT (back propagation through time)

### Truncated Backpropagation through time



# Bidirectional RNN – recorrência bidirecional





## seq2seq model (2014)

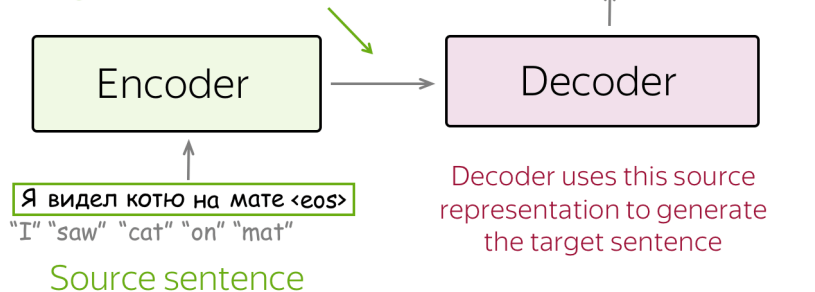
*Sequence to Sequence Learning with Neural Networks*

Ilya Sutskever, Oriol Vinyals, Quoc V. Le

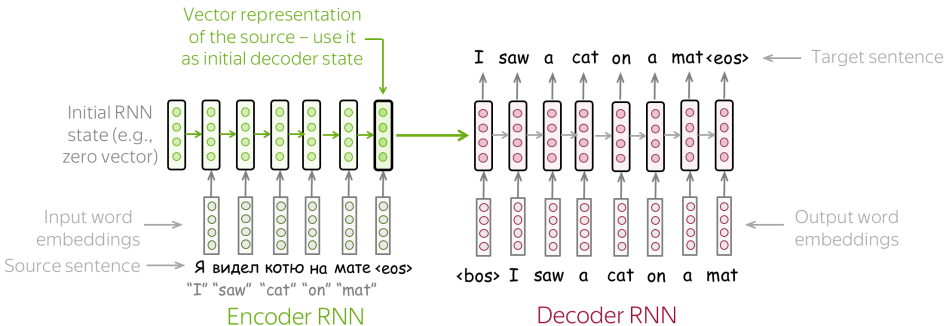
<https://arxiv.org/abs/1409.3215>

- arquitetura encoder-decoder de RNNs

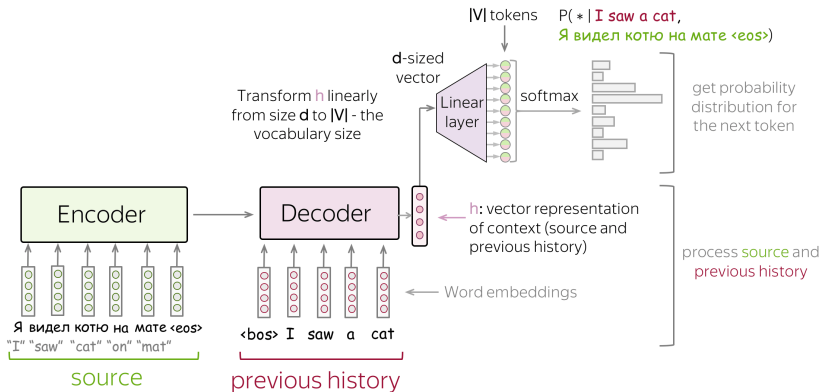
Encoder builds a representation of the source and gives it to the decoder



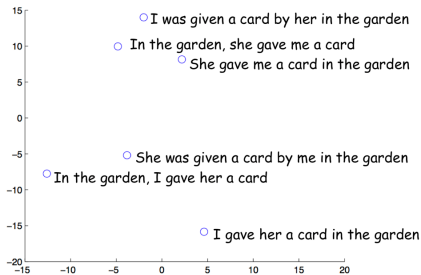
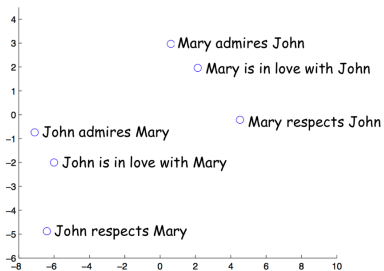
[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)



[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)



[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)



[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)

## seq2seq model with attention (2014)

*Neural Machine Translation by Jointly Learning to Align and Translate*

Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio

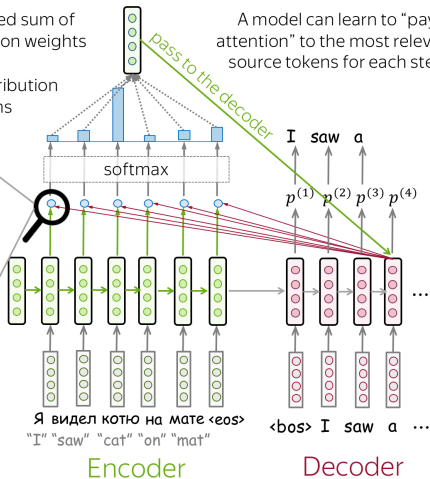
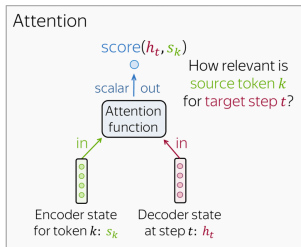
<https://arxiv.org/abs/1409.0473>

- trata o problema de dependências longas

Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

A model can learn to “pay attention” to the most relevant source tokens for each step



[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)

