

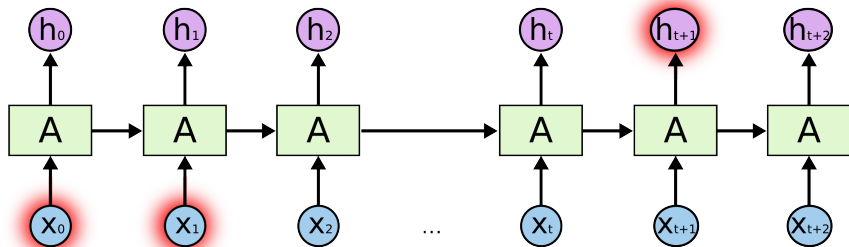
Deep Learning

Vazgen Mikayelyan

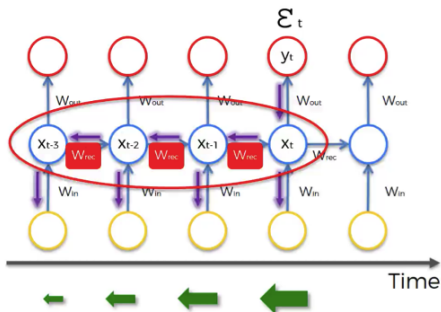
December 1, 2020



Problem of Long Term Dependencies



The Vanishing Gradient Problem



$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\partial \mathcal{E}_t}{\partial \theta} \quad (3)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left(\frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right) \quad (4)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} \mathbf{W}_{rec}^T \text{diag}(\sigma'(\mathbf{x}_{i-1})) \quad (5)$$

$W_{rec} \sim \text{small} \Rightarrow \text{Vanishing}$
 $W_{rec} \sim \text{large} \Rightarrow \text{Exploding}$

Formula Source: Razvan Pascanu et al. (2013)

- 1 GRU and LSTM
- 2 Bidirectional and Deep RNNs
- 3 Attention Models

Simple RNN

Simple RNN

Additional state

Why to use tanh?

- 1 GRU and LSTM
- 2 Bidirectional and Deep RNNs
- 3 Attention Models

Deep RNNs

- 1 GRU and LSTM
- 2 Bidirectional and Deep RNNs
- 3 Attention Models

Encode each word in the sentence into a vector using RNNs.

Encode each word in the sentence into a vector using RNNs.
When decoding, perform a convex combination of these vectors, weighted by attention weights .

Encode each word in the sentence into a vector using RNNs.

When decoding, perform a convex combination of these vectors, weighted by “attention weights”.

Use this combination in picking the next word.

Attention Model

