

Deep Learning and Temporal Data Processing

3 - Recurrent Neural Networks

Andrea Palazzi

June 21th, 2017

University of Modena and Reggio Emilia

Introduction

Vanilla RNN

Training a RNN

Advanced RNN Architectures

Credits

Introduction

In **feedforward neural network** computation flows directly from input \mathbf{x} through intermediate layers \mathbf{h} to output \mathbf{y} .

Conversely, some networks topology feature feedback connections, in other words model outputs are fed back into the model itself.

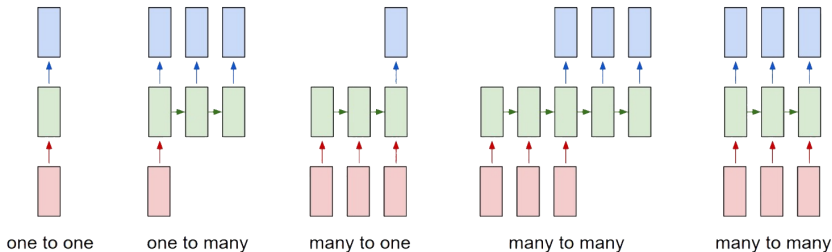
The term **recurrent neural networks** defines this family of models.

Recurrent neural networks (RNN) are **specialized for processing sequences**.

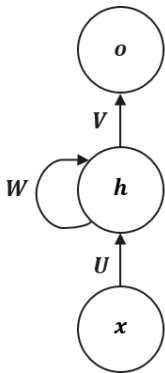
Similarly, we saw that convolutional neural networks feature specialized architecture for processing images.

RNNs boast a **much wider API with respect to feedforward neural networks**.

Indeed, these models can deal with *sequences* in the input, in the output or even both.



Vanilla RNN

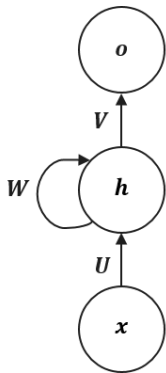


The vanilla RNN is provided with three sets of parameters:

- U maps inputs to the hidden state
- W parametrizes hidden state transition
- V maps hidden state to output

System dynamics is as simple as:

$$\begin{cases} \mathbf{h}^{(t)} = \phi(\mathbf{W} \mathbf{h}^{(t-1)} + \mathbf{U} \mathbf{x}^{(t)}) \\ \mathbf{o}^{(t)} = \mathbf{V} \mathbf{h}^{(t)} \end{cases} \quad (1)$$



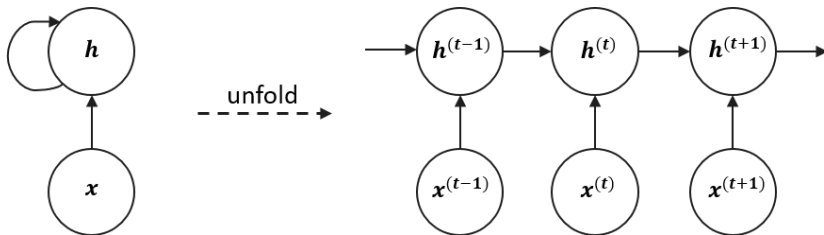
The hidden state $\mathbf{h}^{(t)}$ can be intuitively viewed as a *lossy* summary of the sequence of past inputs fed to the network, in which are stored the main task-relevant aspects of the past sequence of inputs up to time t .

Since the an input sequence of arbitrary length $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)})$ is mapped into a fixed size vector $\mathbf{h}^{(t)}$, this summary is necessarily lossy.

Training a RNN

A recurrent computational graph can be unfolded into a sequential computational graph with a repetitive structure.

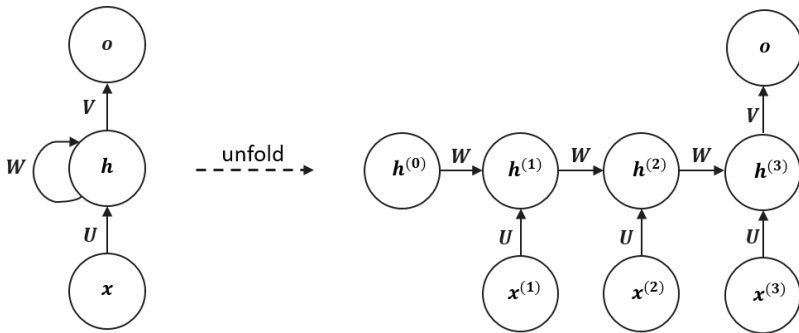
$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta)$$



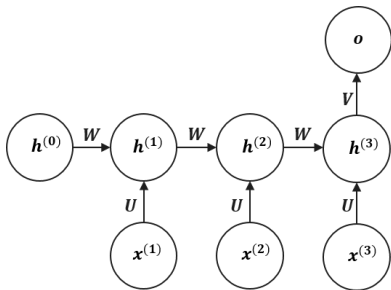
The capacity to **unfold a recurrent graph into a DAG** (Directed Acyclic Graph) allows to train recurrent neural network by means of standard backpropagation.

Because the gradient conceptually flows backward though time instead of though layers, this algorithm is usually referred to as **backpropagation through time (BPTT)**.

Let's make an example for an simple architecture which processes sequences of length $\tau = 3$ and produces an output at the end of the sequence.



Now, given a differentiable loss on final output $L(\mathbf{y}, \mathbf{o})$ let's compute the derivative of the objective L with respect to the weights \mathbf{V} , \mathbf{W} and \mathbf{U} and see what happens.



- $\frac{\partial L}{\partial \mathbf{V}} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \mathbf{V}}$
- $\frac{\partial L}{\partial \mathbf{W}} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \mathbf{h}^{(3)}} \sum_{k=0}^3 \left(\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}} \right)$
- $\frac{\partial L}{\partial \mathbf{U}} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \mathbf{h}^{(3)}} \sum_{k=0}^3 \left(\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{U}} \right)$

Consideration: while $\frac{\partial L}{\partial \mathbf{V}}$ depends only on current state, $\frac{\partial L}{\partial \mathbf{W}}$ and $\frac{\partial L}{\partial \mathbf{U}}$ depend on all previous sequence states.

Looking closer, we see that terms $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}}$ must be themselves computed through the chain rule. For example, we can obtain $\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}}$.

Turns out [3] that when using *tanh* and *sigmoid* activations, the 2-norm of these Jacobian matrices is upper bounded by 1 and 1/4 respectively. Thus, we can easily end up multiplying smaller and smaller numbers until the gradients become zero.

This problem is known as **vanishing gradient problem**, and causes serious trouble when trying to learn long-term dependencies in the input sequences, because contributions of "far-away" steps become zero.

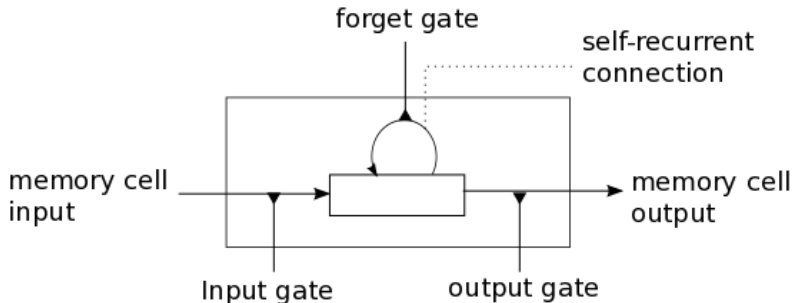
Depending on the network parameters and choice of activation functions, the opposite problem can arise. This is called **exploding gradient problem** and happens when gradients become bigger and bigger until numerical problems destroy the optimization process.

Both issues can be mitigated through proper weight initialization, accurate choice of activation functions and gradient clipping.

Note: these problems can also happen in deep feedforward networks: however, they are more common in recurrent architectures because these models are usually very deep (actually as deep as the length of the input sequence).

Advanced RNN Architectures

Long Short-Term Memory (LSTM) and **Gated Recurrent Unit (GRU)** are more complex recurrent architectures that have been proposed [2, 1] to overcome the issues in the gradient flow and to ease the learning of long-term dependencies thanks to the introduction of **learnable gating mechanisms**.



Let's see **how update equations look like for a LSTM model**. Please notice that LSTM framework, notation is usually slightly different than form vanilla RNN.

$$\left\{ \begin{array}{l} \mathbf{i} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_i + \mathbf{s}^{(t-1)} \mathbf{W}_i) \\ \mathbf{f} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_f + \mathbf{s}^{(t-1)} \mathbf{W}_f) \\ \mathbf{o} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_o + \mathbf{s}^{(t-1)} \mathbf{W}_o) \\ \mathbf{g} = \tanh(\mathbf{x}^{(t)} \mathbf{U}_g + \mathbf{s}^{(t-1)} \mathbf{W}_g) \\ \mathbf{c}^{(t)} = \mathbf{c}^{(t-1)} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i} \\ \mathbf{s}^{(t)} = \tanh(\mathbf{c}^{(t)}) \odot \mathbf{o} \end{array} \right. \quad (2)$$

Here \odot denotes element-wise multiplication. If this sounds complicated, in the following slides we'll parse each of these equations to make a sense out of it.

$$\begin{cases} \mathbf{i} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_i + \mathbf{s}^{(t-1)} \mathbf{W}_i) \\ \mathbf{f} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_f + \mathbf{s}^{(t-1)} \mathbf{W}_f) \\ \mathbf{o} = \sigma(\mathbf{x}^{(t)} \mathbf{U}_o + \mathbf{s}^{(t-1)} \mathbf{W}_o) \\ \dots \end{cases} \quad (3)$$

These are the input, forget and output **gates** respectively. Each gate has the same dimension of the hidden state. Gates are multiplied element-wise with other functions of LSTM thus acting as **continuous, differentiable switches** thanks to the sigmoid activation function.

Input and *output gates* controls how much of the input let through and how much of the internal state exposing to the external, respectively. Conversely, *forget gate* controls how much memory form previous time step must be overwritten.

$$\begin{cases} \dots \\ \mathbf{g} = \tanh(\mathbf{x}^{(t)} \mathbf{U}_g + \mathbf{s}^{(t-1)} \mathbf{W}_g) \\ \dots \end{cases} \quad (4)$$

This equation computes what could intuitively be described a **candidate state**. Indeed, the equation is pretty much the same that we saw for vanilla RNN architecture. Nonetheless, the amount of influence of \mathbf{g} on the LSTM memory cell is controlled by input gate \mathbf{i} .

$$\begin{cases} \dots \\ \mathbf{c}^{(t)} = \mathbf{c}^{(t-1)} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i} \\ \mathbf{s}^{(t)} = \tanh(\mathbf{c}^{(t)}) \odot \mathbf{o} \end{cases} \quad (5)$$

First equation computes the **update for memory cell \mathbf{c}** . The *forget gate* controls how much of memory from previous steps $\mathbf{c}^{(t)}$ must be kept. *Input gate* supervises the amount of newly computed state \mathbf{g} that has to flow into the memory.

Eventually, last equation compute the **output hidden state** from the current memory. *Output gate* regulates the amount of information to be exposed to successive layers.

Credits

These slides heavily borrow from a number of awesome sources. I'm really grateful to all the people who take the time to share their knowledge on this subject with others.

In particular:

- Stanford CS231n Convolutional Neural Networks for Visual Recognition
<http://cs231n.stanford.edu/>
- Deep Learning Book (GoodFellow, Bengio, Courville)
<http://www.deeplearningbook.org/>
- Convolution arithmetic animations
https://github.com/vdumoulin/conv_arithmetic

- Andrej Karphathy personal blog
<http://karpathy.github.io/>
- WildML blog on AI, DL and NLP
<http://www.wildml.com/>
- Michael Nielsen Deep Learning online book
<http://neuralnetworksanddeeplearning.com/>

- [1] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio.

Learning phrase representations using rnn encoder-decoder for statistical machine translation.

arXiv preprint arXiv:1406.1078, 2014.

- [2] S. Hochreiter and J. Schmidhuber.

Long short-term memory.

Neural computation, 9(8):1735–1780, 1997.

[3] R. Pascanu, T. Mikolov, and Y. Bengio.

On the difficulty of training recurrent neural networks.

In *International Conference on Machine Learning*, pages 1310–1318, 2013.