

Deep Learning

Recurrent Neural Network (RNNs)

Ali Ghodsi

University of Waterloo

October 23, 2015

Slides are partially based on Book in preparation, Deep Learning

by Bengio, Goodfellow, and Aaron Courville, 2015

Sequential data

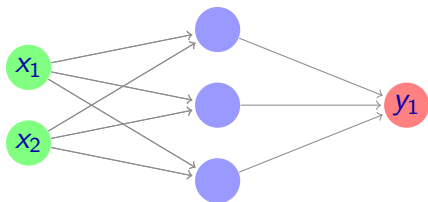
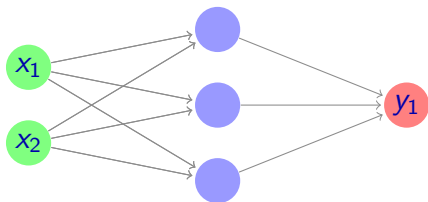
Recurrent neural networks (RNNs) are often used for handling sequential data.

They introduced first in 1986 (Rumelhart et al 1986).

Sequential data usually involves variable length inputs.

Parameter sharing

Parameter sharing makes it possible to extend and apply the model to examples of different lengths and generalize across them.



Recurrent neural network

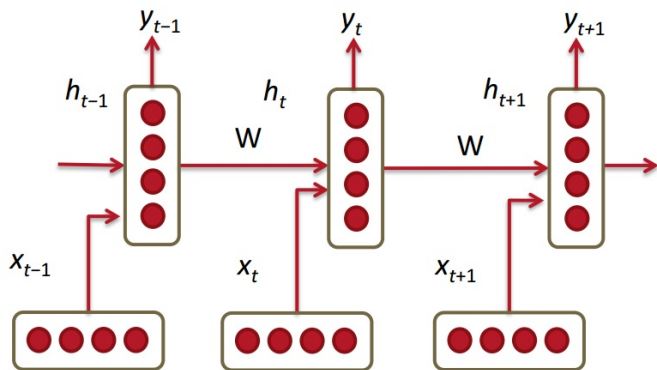
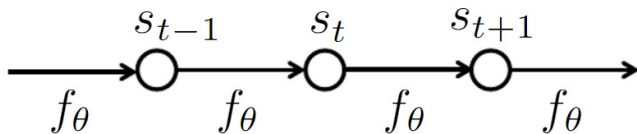


Figure: Richard Socher

Dynamic systems

The classical form of a dynamical system:

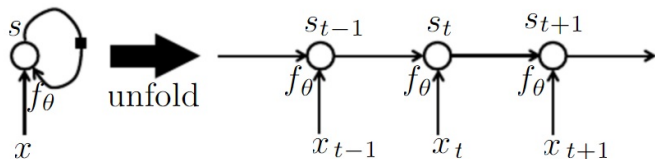
$$s_t = f_{\theta}(s_{t-1})$$



Dynamic systems

Now consider a dynamic system with an external signal x

$$s_t = f_\theta(s_{t-1}, x_t)$$



The state contains information about the whole past sequence.

$$s_t = g_t(x_t, x_{t-1}, x_{t-s}, \dots, x_2, x_1)$$

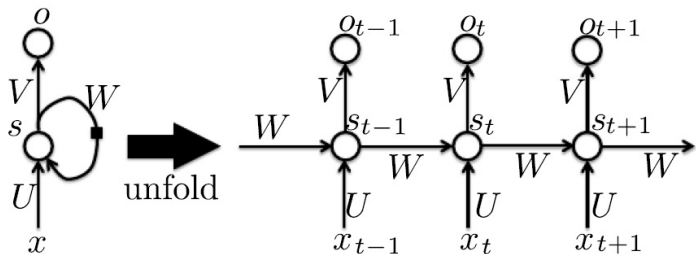
Parameter sharing

We can think of s_t as a summary of the past sequence of inputs up to t .

If we define a different function g_t for each possible sequence length, we would not get any generalization.

If the same parameters are used for any sequence length allowing much better generalization properties.

Recurrent Neural Networks



$$\mathbf{a}_t = \mathbf{b} + W\mathbf{s}_{t-1} + U\mathbf{x}_t$$

$$\mathbf{s}_t = \tanh(\mathbf{a}_t)$$

$$\mathbf{o}_t = \mathbf{c} + V\mathbf{s}_t$$

$$\mathbf{p}_t = \text{softmax}(\mathbf{o}_t)$$

Computing the Gradient in a Recurrent Neural Network

Using the generalized back-propagation algorithm one can obtain the so-called Back-Propagation Through Time (BPTT) algorithm.

Exploding or Vanishing Product of Jacobians

In recurrent nets (also in very deep nets), the final output is the composition of a large number of non-linear transformations.

Even if each of these non-linear transformations is smooth. Their composition might not be.

The derivatives through the whole composition will tend to be either very small or very large.

Exploding or Vanishing Product of Jacobians

The Jacobian (matrix of derivatives) of a composition is the product of the Jacobians of each stage.

If

$$f = f_T \circ f_{T-1} \circ \dots \circ f_2 \circ f_1$$

where $(f \circ g)(x) = f(g(x))$

$$(f \circ g)'(x) = (f' \circ g)(x) \cdot g'(x) = f'(g(x))g'(x)$$

Exploding or Vanishing Product of Jacobians

The Jacobian matrix of derivatives of $f(\mathbf{x})$ with respect to its input vector \mathbf{x} is

$$f' = f'_T f'_{T-1} \dots, f'_2 f_1$$

where

$$f' = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

and

$$f'_t = \frac{\partial f_t(a_t)}{\partial a_t},$$

where $a_t = f_{t-1}(f_{t-2}(\dots, f_2(f_1(\mathbf{x}))))$.

Exploding or Vanishing Product of Jacobians

Simple example

Suppose: all the numbers in the product are scalar and have the same value α .

multiplying many numbers together tends to be either very large or very small.

If T goes to ∞ , then

α^T goes to ∞ if $\alpha > 1$

α^T goes to 0 if $\alpha < 1$

Difficulty of Learning Long-Term Dependencies

Consider a general dynamic system

$$s_t = f_\theta(s_{t-1}, x_t)$$

$$o_t = g_w(s_t),$$

A loss L_t is computed at time step t as a function of o_t and some target y_t . At time T :

$$\frac{\partial L_T}{\partial \theta} = \sum_{t \leq T} \frac{\partial L_t}{\partial s_t} \frac{\partial s_t}{\partial \theta}$$

$$\frac{\partial L_T}{\partial \theta} = \sum_{t \leq T} \frac{\partial L_t}{\partial s_T} \frac{\partial s_T}{\partial s_t} \frac{\partial f_\theta(s_{t-1}, x_t)}{\partial \theta}$$

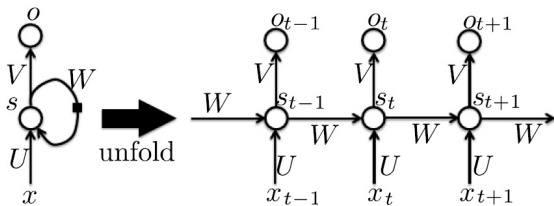
$$\frac{\partial s_T}{\partial s_t} = \frac{\partial s_T}{\partial s_{T-1}} \frac{\partial s_{T-1}}{\partial s_{T-2}} \cdots \frac{\partial s_{t+1}}{\partial s_t}$$

Facing the challenge

Gradients propagated over many stages tend to either vanish (most of the time) or explode.

Echo State Networks

set the recurrent and input weights such that the recurrent hidden units do a good job of capturing the history of past inputs, and only learn the output weights.



$$\mathbf{s}_t = \sigma(W\mathbf{s}_{t-1} + U\mathbf{x}_t)$$

Echo State Networks

If a change Δs in the state at t is aligned with an eigenvector v of jacobian J with eigenvalue $\lambda > 1$, then the small change Δs becomes $\lambda \Delta s$ after one time step, and $\lambda^t \Delta s$ after t time steps.

If the largest eigenvalue $\lambda < 1$, the map from t to $t + 1$ is contractive.

The network forgetting information about the long-term past.

Set the weights to make the Jacobians *slightly contractive*.

Long delays

Use recurrent connections with long delays.

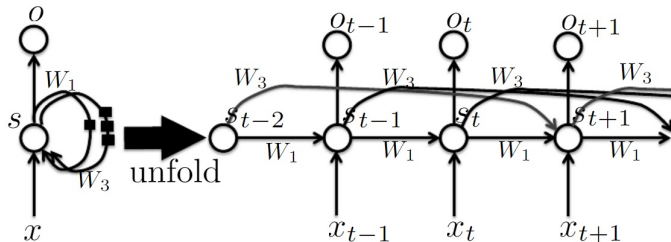


Figure from Benjio et al 2015

Leaky Units

Recall that

$$\mathbf{s}_t = \sigma(W\mathbf{s}_{t-1} + U\mathbf{x}_t)$$

Consider

$$\mathbf{s}_{t,i} = \left(1 - \frac{1}{\tau_i}\right)\mathbf{s}_{t-1} + \frac{1}{\tau_i}\sigma(W\mathbf{s}_{t-1} + U\mathbf{x}_t)$$

$$1 \leq \tau_i \leq \infty$$

$\tau_i = 1$, Ordinary RNN

$\tau_i > 1$, gradients propagate more easily.

$\tau_i \gg 1$, the state changes very slowly, integrating the past values associated with the input sequence.

Gated RNNs

It might be useful for the neural network to forget the old state in some cases.

Example: *a a b b b a a a a b a b*

It might be useful to keep the memory of the past.

Example:

Instead of manually deciding when to clear the state, we want the neural network to learn to decide when to do it.

Gated RNNs, the Long-Short-Term-Memory

The Long-Short-Term-Memory (LSTM) algorithm was proposed in 1997 (Hochreiter and Schmidhuber, 1997).

Several variants of the LSTM are found in the literature:

Hochreiter and Schmidhuber 1997

Graves, 2012

Graves et al., 2013

Sutskever et al., 2014

the principle is always to have a linear self-loop through which gradients can flow for long duration.

Gated Recurrent Units (GRU)

Recent work on gated RNNs, Gated Recurrent Units (GRU) was proposed in 2014

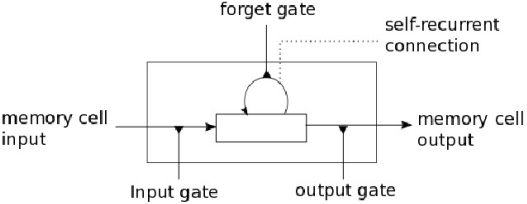
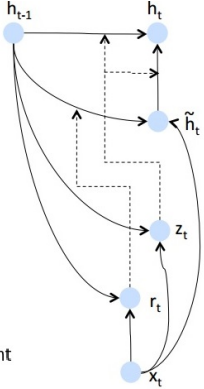
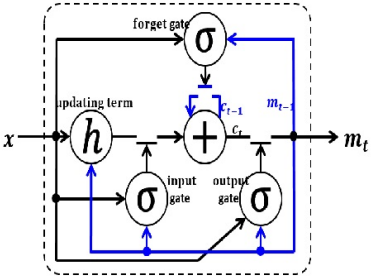
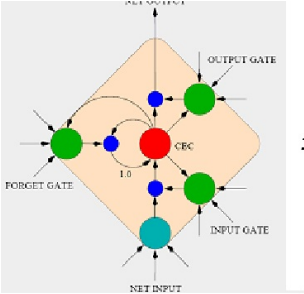
Cho et al., 2014

Chung et al., 2014, 2015

Jozefowicz et al., 2015

Chrupala et al., 2015

Gated RNNs



Gated Recurrent Units (GRU)

Standard RNN computes hidden layer at next time step directly:

$$h_t = f(Wh_{t-1} + Ux_t)$$

GRU first computes an update Gate (another layer) based on current input vector and hidden state

$$z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1})$$

compute reset gate similarly but with different weights

$$r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1})$$

Credit: Richard Socher

Gated Recurrent Units (GRU)

Update gate: $z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1})$

Reset gate: $r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1})$

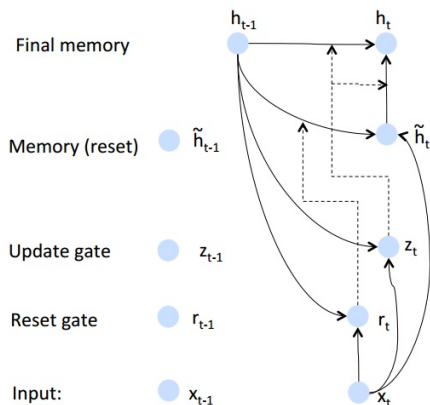
New memory content: $\tilde{h}_t = \tanh(wx_t + r_t \circ Uh_{t-1})$

If reset gate is **0**, then this ignores previous memory and only stores the new information

Final memory at time step combines current and previous time steps:

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

Gated Recurrent Units (GRU)



Richard Socher

Update gate:

$$z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1})$$

Reset gate:

$$r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1})$$

New memory content:

$$\tilde{h}_t = \tanh(wx_t + r_t \circ Uh_{t-1})$$

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

Gated Recurrent Units (GRU)

$$\begin{aligned}z_t &= \sigma(W^{(z)}x_t + U^{(z)}h_{t-1}) \\r_t &= \sigma(W^{(r)}x_t + U^{(r)}h_{t-1}) \\\tilde{h}_t &= \tanh(Wx_t + r_t \circ Uh_{t-1}) \\h_t &= z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t\end{aligned}$$

If reset is close to 0, ignore previous hidden state \rightarrow Allow model to drop information that is irrelevant

Update gate z controls how much of past state should matter now. If z close to 1, then we can copy information in that unit through many time steps.

Units with short term dependencies often have reset gates very active.

The Long-Short-Term-Memory (LSTM)

We can make the units even more complex

Allow each time step to modify

Input gate (current cell matters) $i_t = \sigma(W^{(i)}x_t + U^{(i)}h_{t-1})$

Forget (gate 0, forget past) $f_t = \sigma(W^{(f)}x_t + U^{(f)}h_{t-1})$

Output (how much cell is exposed) $o_t = \sigma(W^{(o)}x_t + U^{(o)}h_{t-1})$

New memory cell $\tilde{c}_t = \tanh(W^{(c)}x_t + U^{(c)}h_{t-1})$

Final memory cell: $c_t = f_t \circ c_{t-1} + (i_t) \circ \tilde{c}_t$

Final hidden state: $h_t = o_t \circ \tanh(c_t)$

Clipping Gradients

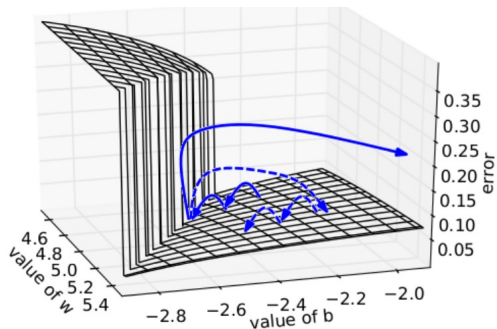


Figure: Pascanu et al., 2013

Strongly non-linear functions tend to have derivatives that can be either very large or very small in magnitude.

Clipping Gradients

Simple solution for clipping the gradient. (Mikolov, 2012; Pascanu et al., 2013):

Clip the parameter gradient from a mini batch element-wise (Mikolov, 2012) just before the parameter update.

Clip the norm $\|g\|$ of the gradient g (Pascanu et al., 2013a) just before the parameter update.