

Long Short-Term Memory (LSTM)

A brief introduction

Daniel Renshaw

24th November 2014

(updated 30th November 2015)

Introduction

- ▶ What is Long Short-Term Memory (LSTM)?
 - ▶ A type of artificial neural network (ANN) layer
 - ▶ More specifically: a type of *recurrent* ANN layer
- ▶ Not a whole network architecture, just a layer
- ▶ Designed to solve the vanishing/exploding gradient problem
- ▶ *Why does it look so complex?*

Context and notation

- ▶ Task: ANN language modelling (example for this presentation)
- ▶ Word vocabulary, size V
- ▶ $\mathbf{x}_t \in \mathbb{R}^V$: true word in position t (one-hot)
- ▶ $\mathbf{y}_t \in \Delta^{V-1}$: predicted word in position t (distribution)
- ▶ Assume inputs are sentences zero padded to length L

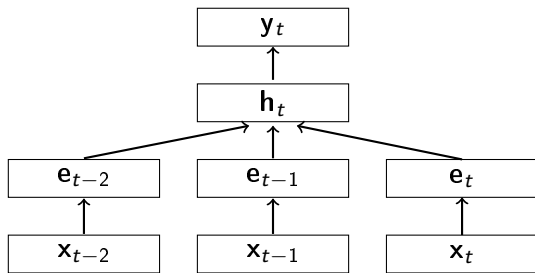
Context and notation

- ▶ Model: $\mathbf{y}_t = p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \dots, \mathbf{x}_1)$ for $1 \leq t < L$
- ▶ Minimize cross-entropy objective:

$$J = \sum_{t=1}^{L-1} H(\mathbf{y}_t, \mathbf{x}_{t+1}) = \sum_{t=1}^{L-1} \sum_{i=1}^V \log(y_{t,i}) x_{t+1,i}$$

- ▶ $\sigma(\bullet)$ is a sigmoid-shaped function (e.g. *logistic* or *tanh*)
- ▶ \mathbf{b}^\bullet is a bias vector, W^\bullet is a weight matrix

Multi-Layer Perceptron (MLP)

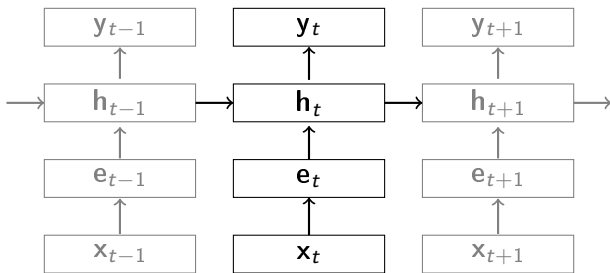


$$\mathbf{y}_t = \text{softmax} \left(W^{yh} \mathbf{h}_t + \mathbf{b}^y \right)$$

$$\mathbf{h}_t = \sigma \left(W^{he} [\mathbf{e}_t; \mathbf{e}_{t-1}; \mathbf{e}_{t-2}] + \mathbf{b}^h \right)$$

$$\mathbf{e}_t = W^{\text{ex}} \mathbf{x}_t$$

Recurrent Neural Network (RNN)



$$y_t = \text{softmax} \left(W^{yh} h_t + \mathbf{b}^y \right)$$

$$h_t = \sigma \left(W^{he} e_t + W^{hh} h_{t-1} + \mathbf{b}^h \right)$$

$$e_t = W^{ex} x_t$$

Vanishing gradients

- ▶ Error gradients pass through nonlinearity every step

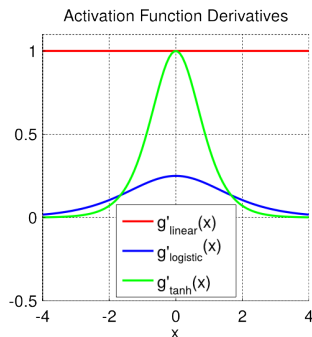
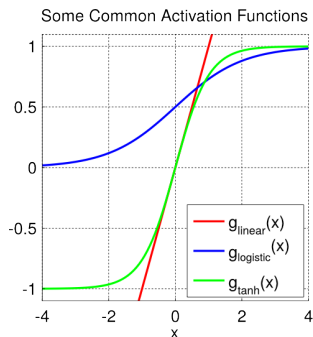


Image from <https://theclevermachine.wordpress.com>

- ▶ Unless weights large, error signal will degrade

$$\delta_h = \sigma'(\bullet) W^{(h+1)h} \delta_{h+1}$$

Vanishing gradients

- ▶ Gradients may vanish or explode
- ▶ Can affect any 'deep' network
 - ▶ e.g. fine-tuning a non-recurrent deep neural network

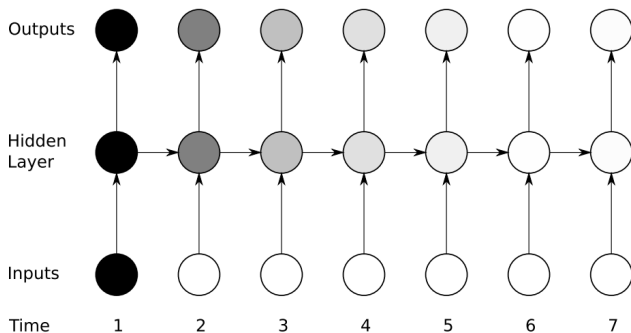
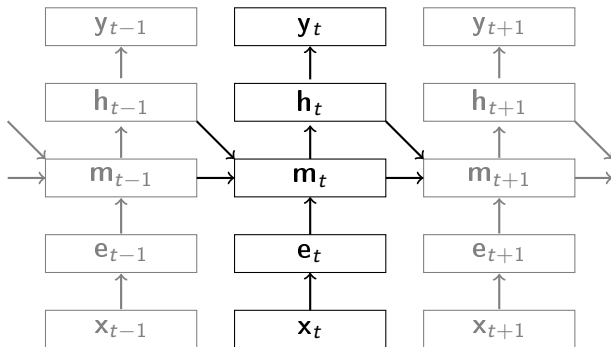


Image from Alex Graves' textbook

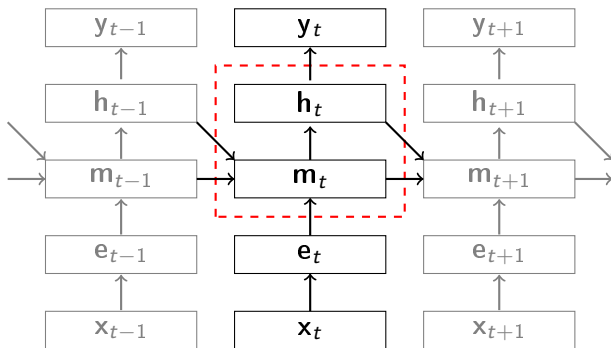
Constant Error Carousel

- ▶ Allow the network to propagate errors without modification
- ▶ No nonlinearity in recursion



Constant Error Carousel

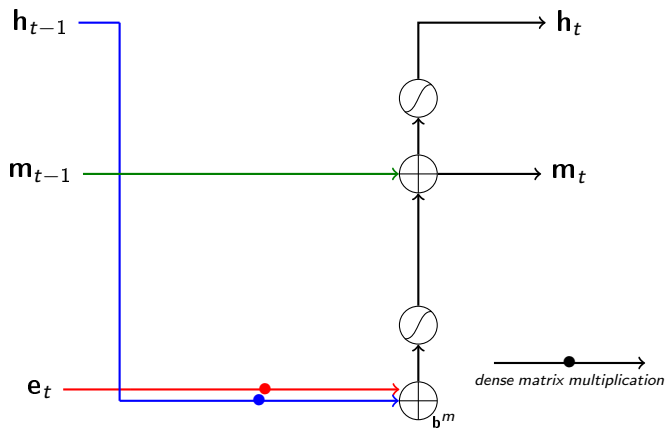
- ▶ Allow the network to propagate errors without modification
- ▶ No nonlinearity in recursion



- ▶ Focus on region in red dashed box from now on

Constant Error Carousel

- ▶ Allow the network to propagate errors without modification
- ▶ No nonlinearity in recursion



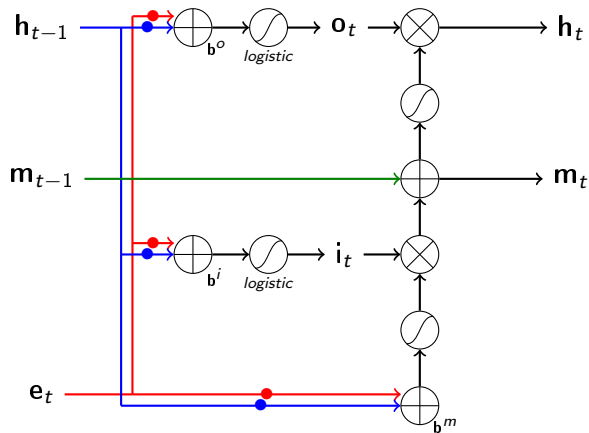
Constant Error Carousel

- ▶ Allow the network to propagate errors without modification
- ▶ No nonlinearity in recursion

$$\mathbf{h}_t = \sigma(\mathbf{m}_t)$$
$$\mathbf{m}_t = \mathbf{m}_{t-1} + \sigma\left(W^{me}\mathbf{e}_t + W^{mh}\mathbf{h}_{t-1} + \mathbf{b}^m\right)$$

LSTM v1: input and output gates

- ▶ Attenuate input and output signals
- ▶ Gates applied elementwise



LSTM v1: input and output gates

- ▶ Attenuate input and output signals
- ▶ Gates applied elementwise

$$\mathbf{h}_t = \mathbf{o}_t \odot \sigma(\mathbf{m}_t)$$

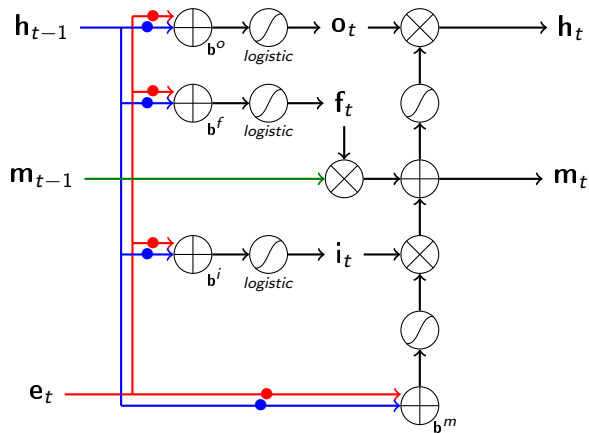
$$\mathbf{o}_t = \text{logistic} \left(W^{oe} \mathbf{e}_t + W^{oh} \mathbf{h}_{t-1} + \mathbf{b}^o \right)$$

$$\mathbf{m}_t = \mathbf{m}_{t-1} + \mathbf{i}_t \odot \sigma \left(W^{me} \mathbf{e}_t + W^{mh} \mathbf{h}_{t-1} + \mathbf{b}^m \right)$$

$$\mathbf{i}_t = \text{logistic} \left(W^{ie} \mathbf{e}_t + W^{ih} \mathbf{h}_{t-1} + \mathbf{b}^i \right)$$

LSTM v2: forget (remember) gate

- ▶ Model controls when memory, \mathbf{m}_t , is reduced
- ▶ Forget gate should be called remember gate



LSTM v2: forget (remember) gate

- ▶ Model controls when memory, \mathbf{m}_t , is reduced
- ▶ Forget gate should be called remember gate

$$\mathbf{h}_t = \mathbf{o}_t \odot \sigma(\mathbf{m}_t)$$

$$\mathbf{o}_t = \text{logistic}\left(W^{oe}\mathbf{e}_t + W^{oh}\mathbf{h}_{t-1} + \mathbf{b}^o\right)$$

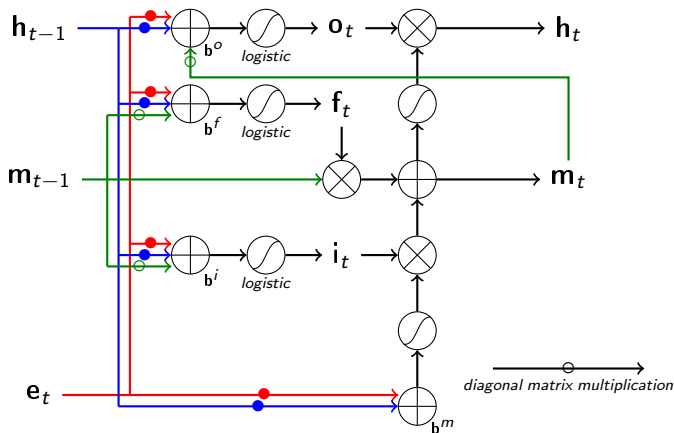
$$\mathbf{m}_t = \mathbf{f}_i \odot \mathbf{m}_{t-1} + \mathbf{i}_t \odot \sigma\left(W^{me}\mathbf{e}_t + W^{mh}\mathbf{h}_{t-1} + \mathbf{b}^m\right)$$

$$\mathbf{i}_t = \text{logistic}\left(W^{ie}\mathbf{e}_t + W^{ih}\mathbf{h}_{t-1} + \mathbf{b}^i\right)$$

$$\mathbf{f}_i = \text{logistic}\left(W^{fe}\mathbf{e}_t + W^{fh}\mathbf{h}_{t-1} + \mathbf{b}^f\right)$$

LSTM v3: peepholes

- ▶ Allow the gates to additionally see the internal memory state
- ▶ Diagonal matrices only (all others dense)



LSTM v3: peepholes

- ▶ Allow the gates to additionally see the internal memory state
- ▶ Diagonal matrices only (all others dense)

$$\mathbf{h}_t = \mathbf{o}_t \odot \sigma(\mathbf{m}_t)$$

$$\mathbf{o}_t = \text{logistic} \left(W^{oe} \mathbf{e}_t + W^{oh} \mathbf{h}_{t-1} + W^{om} \mathbf{m}_t + \mathbf{b}^o \right)$$

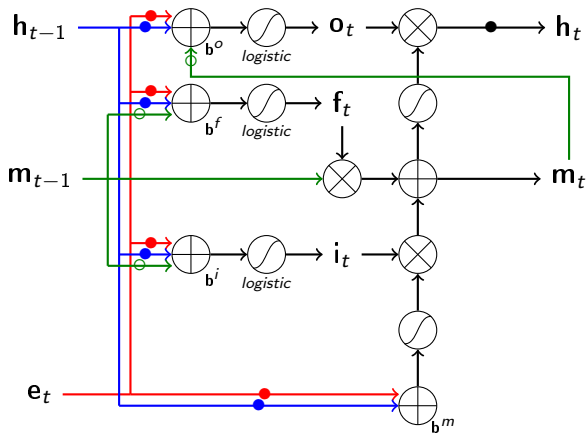
$$\mathbf{m}_t = \mathbf{f}_i \odot \mathbf{m}_{t-1} + \mathbf{i}_t \odot \sigma \left(W^{me} \mathbf{e}_t + W^{mh} \mathbf{h}_{t-1} + \mathbf{b}^m \right)$$

$$\mathbf{i}_t = \text{logistic} \left(W^{ie} \mathbf{e}_t + W^{ih} \mathbf{h}_{t-1} + W^{im} \mathbf{m}_{t-1} + \mathbf{b}^i \right)$$

$$\mathbf{f}_i = \text{logistic} \left(W^{fe} \mathbf{e}_t + W^{fh} \mathbf{h}_{t-1} + W^{fm} \mathbf{m}_{t-1} + \mathbf{b}^f \right)$$

LSTM v4: output projection layer

- ▶ Reduces dimensionality of recursive messages
- ▶ Can speed up training without affecting results quality



LSTM v4: output projection layer

- ▶ Reduces dimensionality of recursive messages
- ▶ Can speed up training without affecting results quality

$$\mathbf{h}_t = W^{hm} (\mathbf{o}_t \odot \sigma(\mathbf{m}_t))$$

$$\mathbf{o}_t = \text{logistic} \left(W^{oe} \mathbf{e}_t + W^{oh} \mathbf{h}_{t-1} + W^{om} \mathbf{m}_t + \mathbf{b}^o \right)$$

$$\mathbf{m}_t = \mathbf{f}_i \odot \mathbf{m}_{t-1} + \mathbf{i}_t \odot \sigma \left(W^{me} \mathbf{e}_t + W^{mh} \mathbf{h}_{t-1} + \mathbf{b}^m \right)$$

$$\mathbf{i}_t = \text{logistic} \left(W^{ie} \mathbf{e}_t + W^{ih} \mathbf{h}_{t-1} + W^{im} \mathbf{m}_{t-1} + \mathbf{b}^i \right)$$

$$\mathbf{f}_i = \text{logistic} \left(W^{fe} \mathbf{e}_t + W^{fh} \mathbf{h}_{t-1} + W^{fm} \mathbf{m}_{t-1} + \mathbf{b}^f \right)$$

Gradients no longer vanish

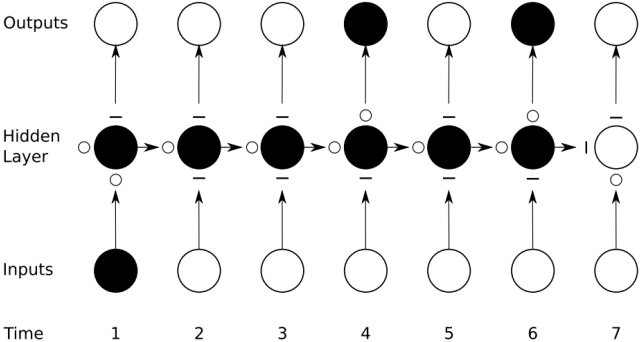


Image from Alex Graves' textbook

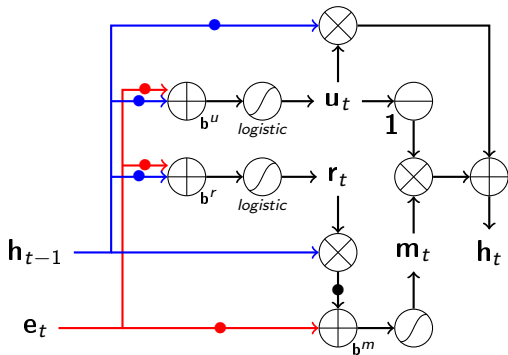
LSTM implementations

- ▶ RNNLIB (Alex Graves) <http://sourceforge.net/p/rnnl/>
- ▶ Built-in: Lasagne, Blocks, Keras
- ▶ Built-out: Caffe (to be built-in later), Torch, Theano, e.g.

```
def lstm_step(x_t, m_tm1, h_tm1, w_xe, ..., b_y):
    e_t = dot(x_t, w_xe)
    i_t = sigmoid(dot(e_t, w_ei) + dot(m_tm1, w_mi) + c_tm1 * w_ci + b_i)
    f_t = sigmoid(dot(e_t, w_ef) + dot(m_tm1, w_mf) + c_tm1 * w_cf + b_f)
    h_t = f_t * h_tm1 + i_t * tanh(dot(e_t, w_eh) + dot(m_tm1, w_mh) + b_h)
    o_t = sigmoid(dot(e_t, w_eo) + dot(m_tm1, w_mo) + c_t * w_co + b_o)
    m_t = dot(o_t * tanh(h_t), w_mm)
    y_t = softmax(dot(m_t, w_my) + b_y)
    return m_t, h_t, y_t
```

Gated Recurrent Unit

- ▶ Introduced with RNN Encoder-Decoder (Cho et al. 2014)
- ▶ Complexity midway between RNN and LSTM



Gated Recurrent Unit

- ▶ Introduced with RNN Encoder-Decoder (Cho et al. 2014)
- ▶ Complexity midway between RNN and LSTM

$$\mathbf{h}_t = \mathbf{u}_t \odot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{u}_t) \mathbf{m}_t$$

$$\mathbf{u}_t = \text{logistic} \left(W^{ue} \mathbf{e}_t + W^{uh} \mathbf{h}_{t-1} + \mathbf{b}^u \right)$$

$$\mathbf{m}_t = \sigma \left(W^{me} \mathbf{e}_t + W^{mh} (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}^m \right)$$

$$\mathbf{r}_t = \text{logistic} \left(W^{re} \mathbf{e}_t + W^{rh} \mathbf{h}_{t-1} + \mathbf{b}^r \right)$$