

Long Short-Term Memory

Managing Long-Term Dependencies within Sequences

August 25, 2013

Outline

Recurrent Neural Networks

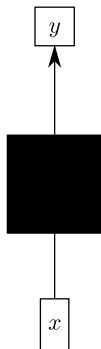
Long Short-Term Memory Cells

Example

Output layer for text recognition

Classification task

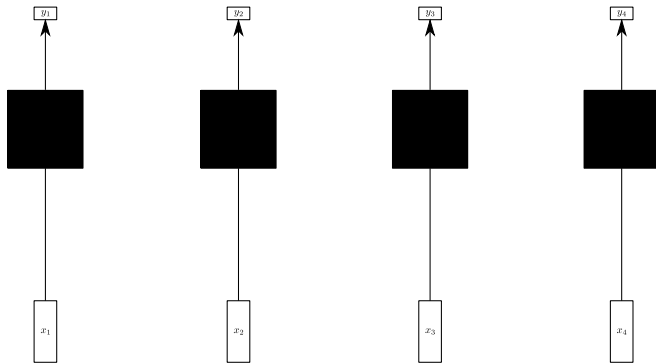
mapping x to y



$$\mathcal{F} : x \in \mathbb{X} \rightarrow y \in \mathbb{Y}$$

Sequence classification task

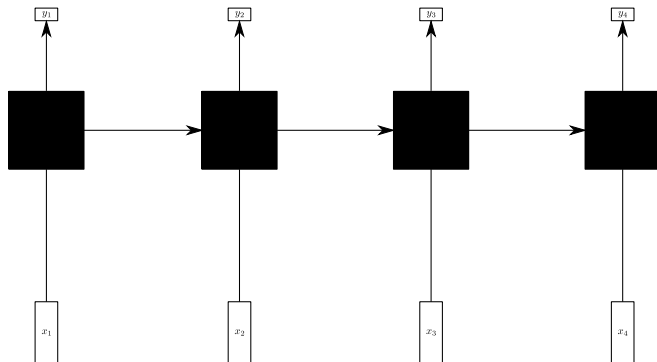
mapping $x_1x_2 \dots$ to $y_1y_2 \dots$



$$\mathcal{F} : \mathbf{x} \in \mathbb{X}^* \rightarrow \mathbf{y} \in \mathbb{Y}^*$$

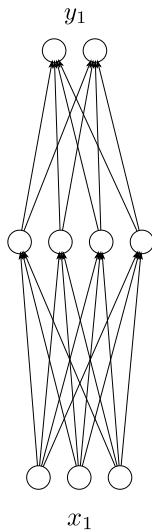
Sequence classification task with dependencies

mapping $x_1x_2 \dots$ to $y_1y_2 \dots$

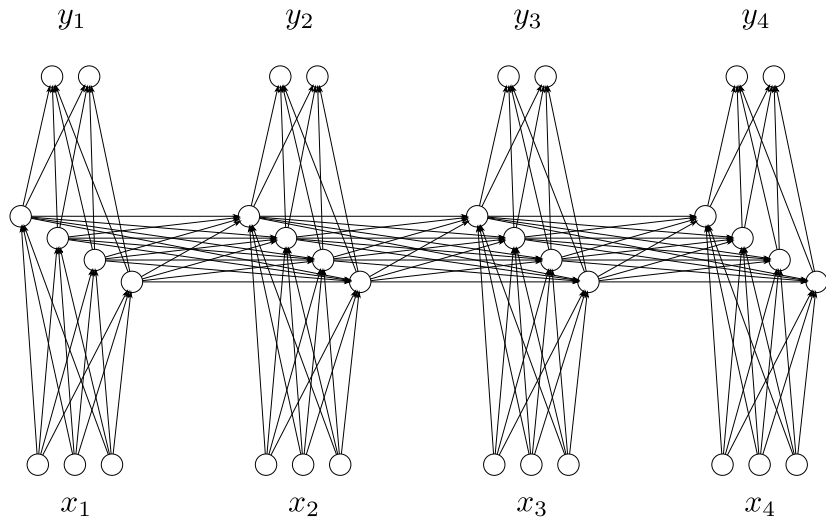


$$\mathcal{F} : \mathbf{x} \in \mathbb{X}^* \rightarrow \mathbf{y} \in \mathbb{Y}^*$$

Feed-forward Neural Networks



Recurrent Neural Networks

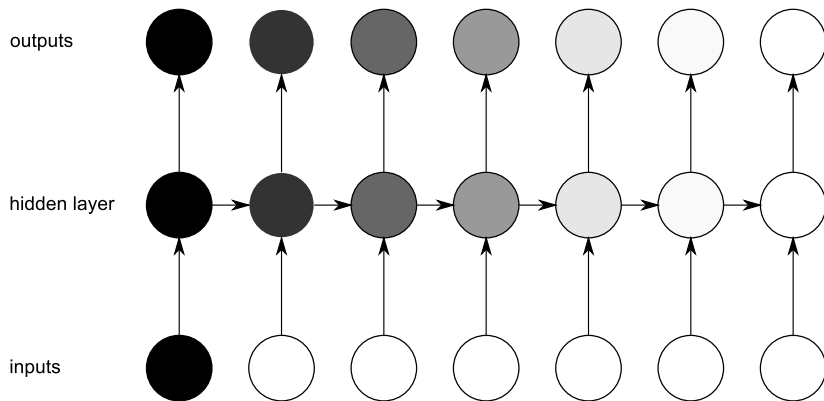


Problems with simple RNN architectures

Problems with simple RNN architectures

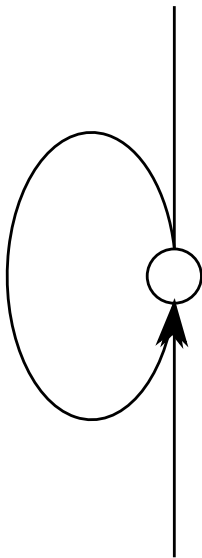
- ▶ Vanishing gradient in training
 - ▶ Sensitivity to an input at time t decreases rapidly
- ▶ Only limited context effectively usable

Vanishing gradient problem



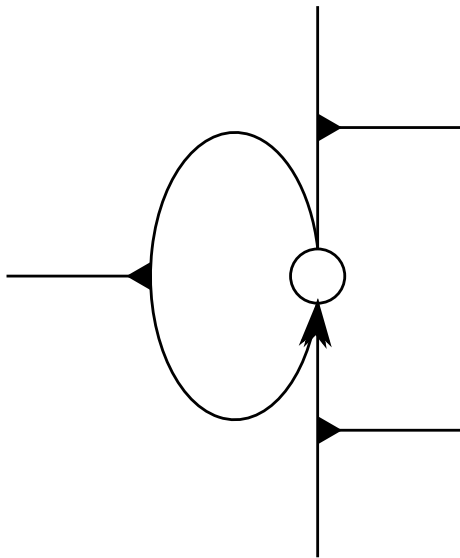
Constructing a differentiable memory cell

Start point



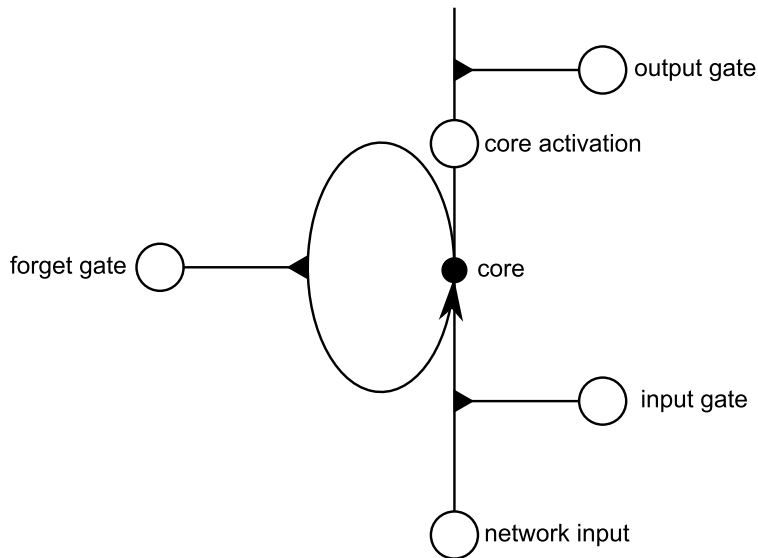
Constructing a differentiable memory cell

Adding controls



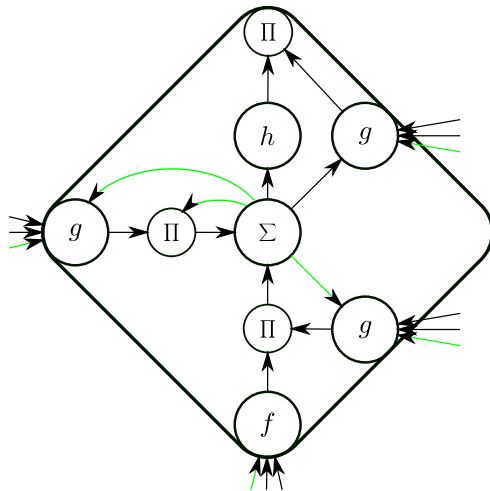
Constructing a differentiable memory cell

Let the controls be neurons



Constructing a differentiable memory cell

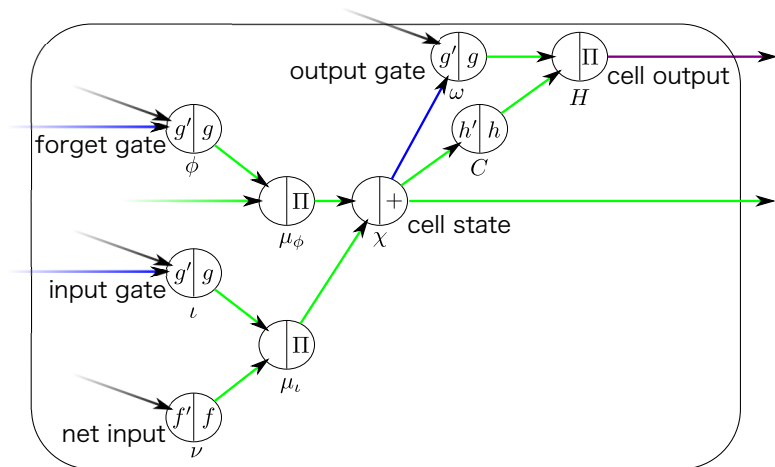
The final result



green lines are time-delayed connections

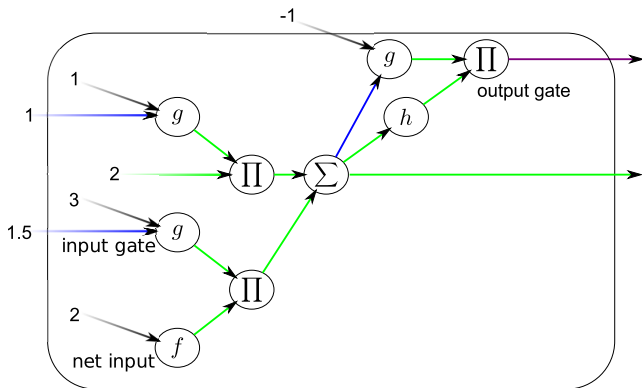
Constructing a differentiable memory cell

The final result



green lines are connection with a fixed weight = 1

Example

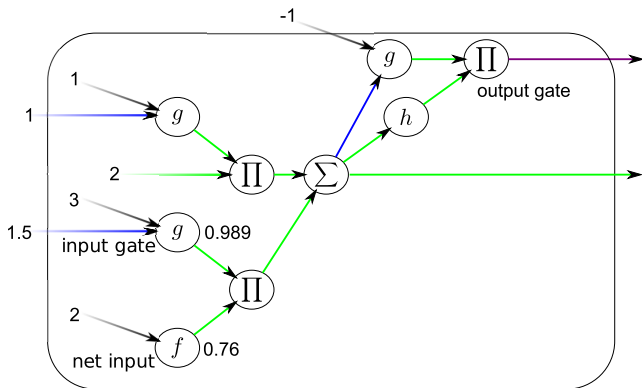


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

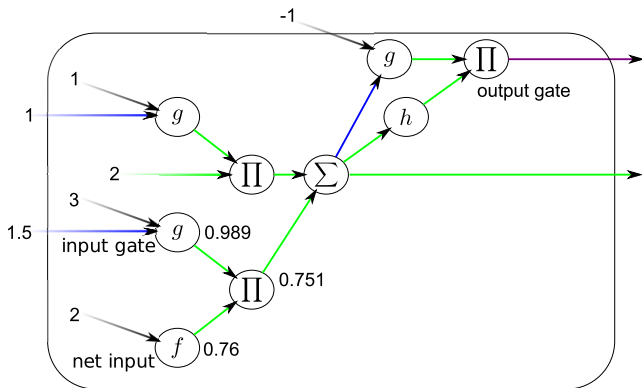


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

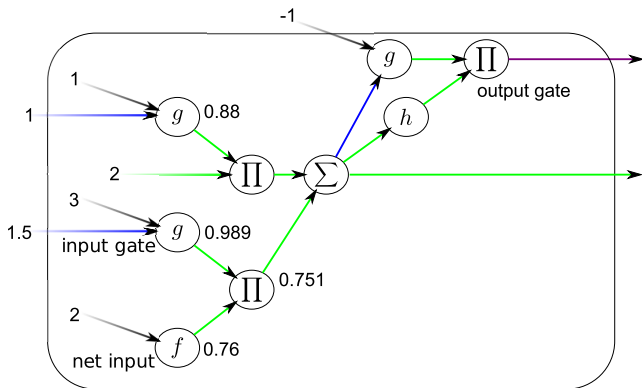


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

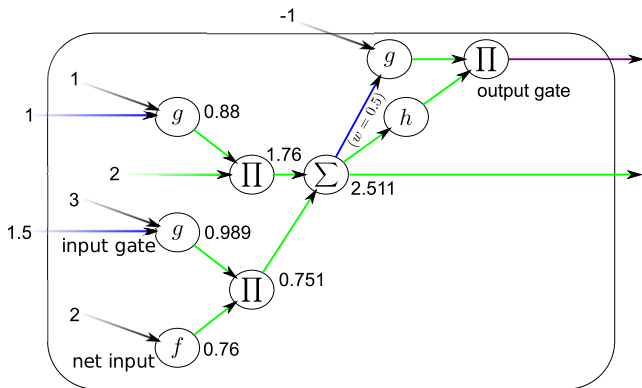


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

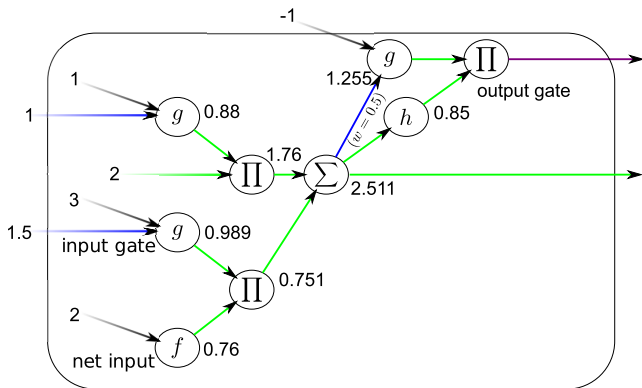


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

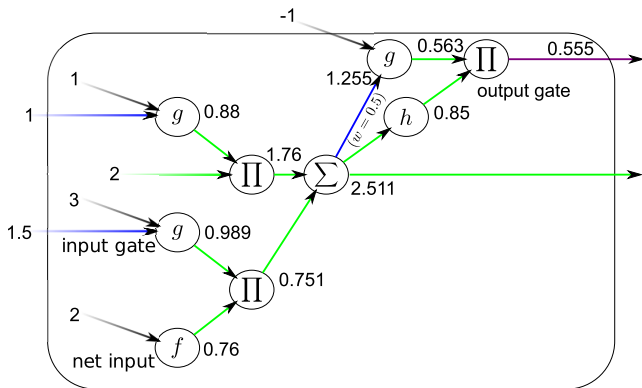


$$g(x) = (1 + \exp(-x))^{-1}$$

$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

Example

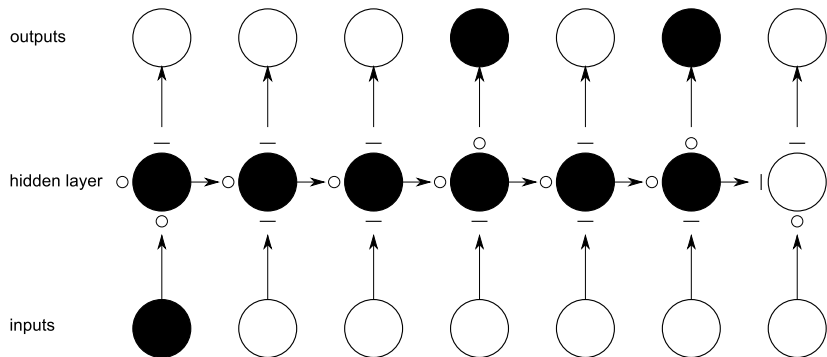


$$g(x) = (1 + \exp(-x))^{-1}$$

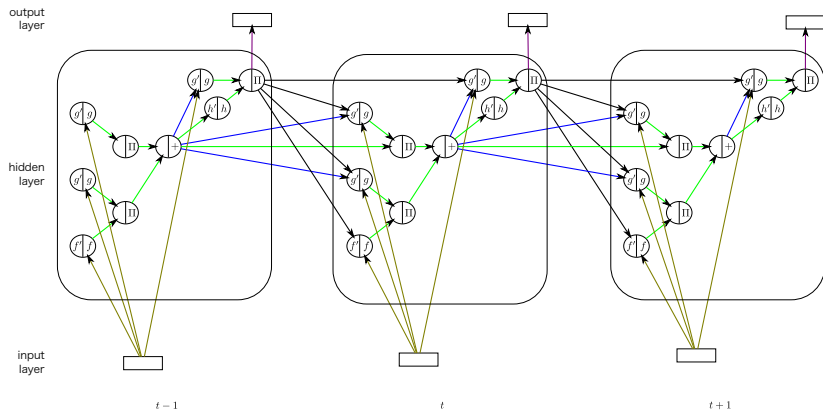
$$f(x) = 2(1 + \exp(-x))^{-1} - 1$$

$$h(x) = 2(1 + \exp(-x))^{-1} - 1$$

LSTM preserving information



Complete Architecture

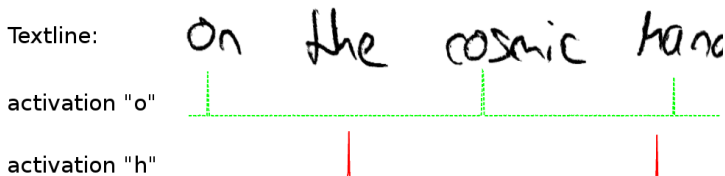


Output Layer

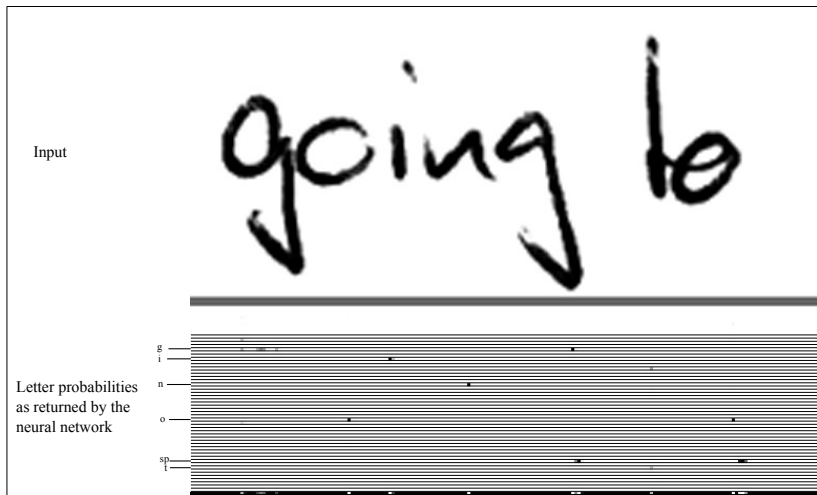
- ▶ The output layer can be any arbitrary output layer
- ▶ For text recognition, one output node is associated to each recognizable label
- ▶ The output layer is usually normalized via *softmax*

$$y_c = \frac{\exp(a_c)}{\sum_{c'} \exp(a_{c'})}$$

- ▶ An extra node, the blank or ε -node, in the output layer can be used as a default output

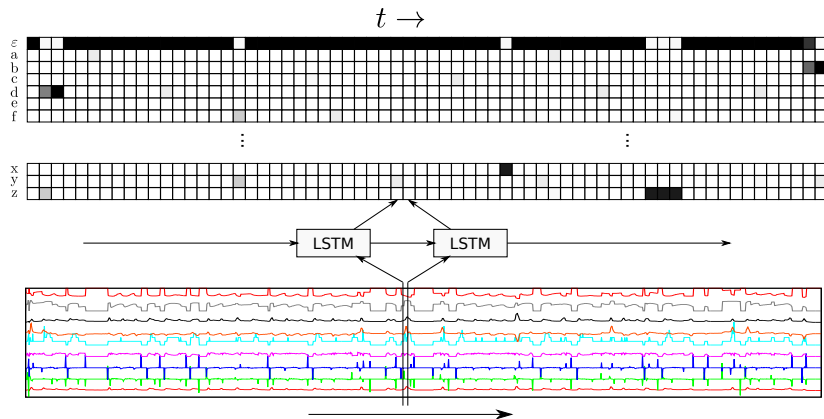


Real Word Example



Complete Architecture

The entire LSTM neural network for text recognition



Training a Neural Network

The process for training a (already initialized) neural network is

- ▶ take a new training example (\mathbf{x}, \mathbf{z})
- ▶ Let the network produce an output $\mathbf{y} = NN(\mathbf{x})$
- ▶ Compute the error $E(\mathbf{y}, \mathbf{z})$
- ▶ Compute the derivative of the error wrt. to every weight
 $\delta w = \frac{\partial E(\mathbf{x}, \mathbf{y})}{\partial w}$
- ▶ Change each weight to reduce the error $w \leftarrow w - \eta \cdot \delta w$

Summary

LSTM

- ▶ Recurrent neural network
- ▶ Gates control the flow of information
- ▶ Differentiable version of a memory cell
- ▶ No vanishing gradient problem
- ▶ Trainable to consider long-term dependencies

Forward Pass

- ▶ Simple addition and multiplication operations
- ▶ Runs in $O(|x| \cdot (|LSTM|^2 + |x| \cdot |LSTM| + |LSTM| \cdot |y|))$

Backward Pass

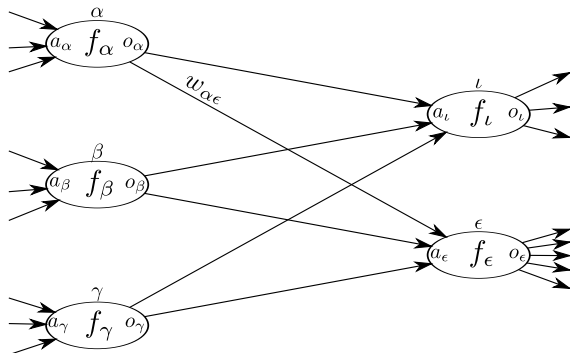
- ▶ Standard Back-propagation
- ▶ Complete unfolding in time feasible
- ▶ Different error functions possible, according to task

Time for Questions

Appendix: Back-propagation Training Formulas

Back-propagation

Error gradient wrt. weight



$$\frac{\partial E}{\partial a} = \delta_a$$

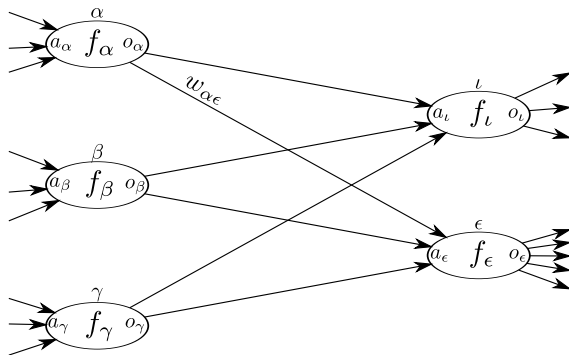
$$\frac{\partial E}{\partial w} = \delta_w$$

$$\frac{\partial E}{\partial o} = \delta_o$$

$$\begin{aligned}\frac{\partial E}{\partial w_{\alpha\epsilon}} &= \frac{\partial E}{\partial a_\epsilon} \frac{\partial a_\epsilon}{\partial w_{\alpha\epsilon}} \\ &= \delta_{a_\epsilon} \cdot o_\alpha\end{aligned}$$

Back-propagation

Error Gradient wrt. node output



$$\frac{\partial E}{\partial a} = \delta_a$$

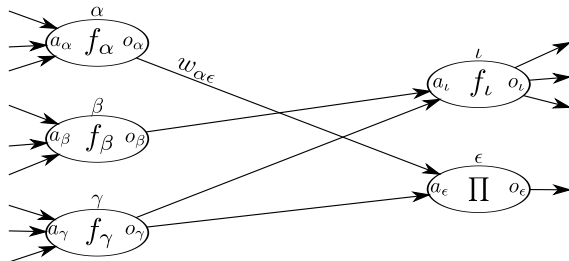
$$\frac{\partial E}{\partial w} = \delta_w$$

$$\frac{\partial E}{\partial o} = \delta_o$$

$$\begin{aligned}\delta_{o_\alpha} = \frac{\partial E}{\partial o_\alpha} &= \sum_k \frac{\partial a_k}{\partial o_\alpha} \frac{\partial E}{\partial a_k} \quad k = \iota, \epsilon \\ &= \sum_k w_{\alpha k} \delta_{a_k} \quad k = \iota, \epsilon\end{aligned}$$

Back-propagation

Error Gradient wrt. node output (for multiplication nodes)



$$\frac{\partial E}{\partial a} = \delta_a$$

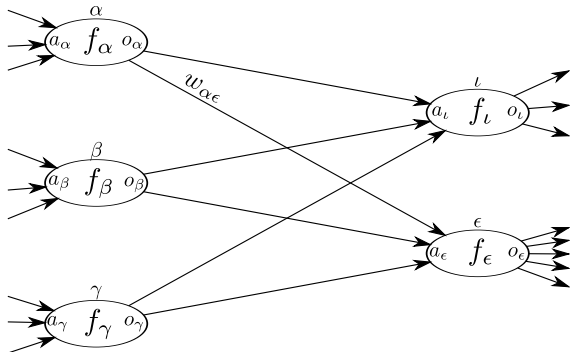
$$\frac{\partial E}{\partial w} = \delta_w$$

$$\frac{\partial E}{\partial o} = \delta_o$$

$$\begin{aligned}\delta_{o_\alpha} &= \frac{\partial E}{\partial o_\alpha} = \frac{\partial E}{\partial o_\epsilon} \frac{\partial o_\epsilon}{\partial o_\alpha} \\ &= \delta_{o_\epsilon} \cdot o_\gamma \cdot w_{\alpha\epsilon} \cdot w_{\gamma\epsilon}\end{aligned}$$

Back-propagation

Error Gradient wrt. node input



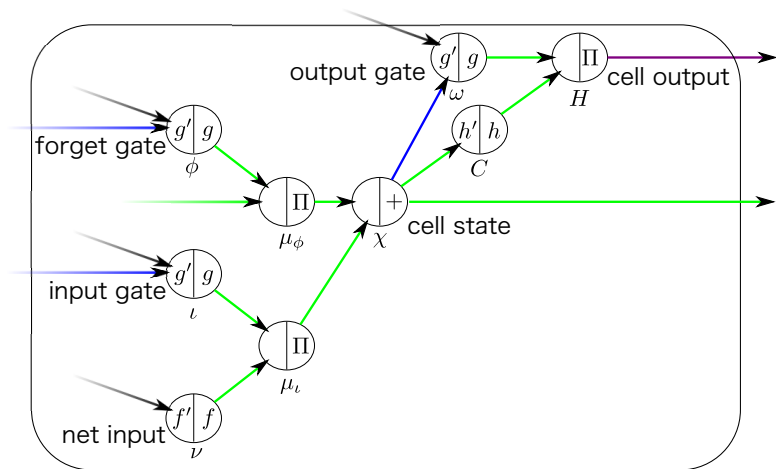
$$\frac{\partial E}{\partial a} = \delta_a$$

$$\frac{\partial E}{\partial w} = \delta_w$$

$$\frac{\partial E}{\partial o} = \delta_o$$

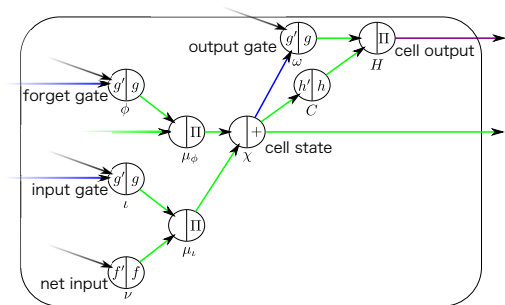
$$\begin{aligned} \frac{\partial E}{\partial a_\alpha} &= \frac{\partial o_\alpha}{\partial a_\alpha} \frac{\partial E}{\partial o_\alpha} \\ &= f'_\alpha(a_\alpha) \sum_k w_{\alpha k} \delta a_k \quad k = l, \epsilon \end{aligned}$$

Deriving the back-propagation formulas for the LSTM cell



Deriving the back-propagation formulas for the LSTM cell

Example: peephole weight output gate

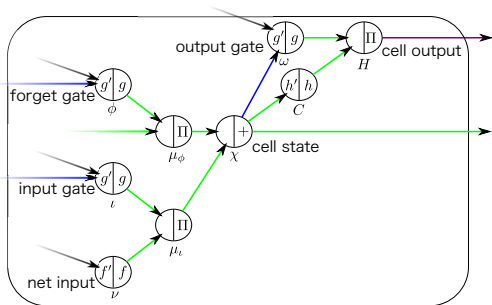


The error gradient at the output of cell output

$$\delta_H^t = \frac{\partial E}{\partial o_{H_i}^t} = \sum_{k \in \text{OutputLayer}} \delta_{a_k} w_{ik}^{(\text{HO})} + \sum_{k \in \text{hiddenLayer}} \delta_{a_k} w_{ih}^{(\text{HH})}$$

Deriving the back-propagation formulas for the LSTM cell

Example: peephole weight output gate

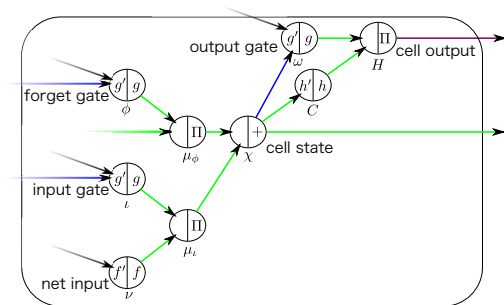


The error gradient
at the input of out-
put gate

$$\delta_\omega = g'(a_\omega^t) \cdot h(a_C^t) \cdot \delta_H^t$$

Deriving the back-propagation formulas for the LSTM cell

Example: peephole weight output gate

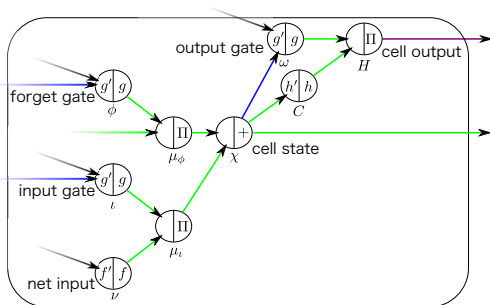


The error gradient
at the input of core
activation

$$\delta_C = h'(o_\chi^t) \cdot o_\omega^t \cdot \delta_H^t$$

Deriving the back-propagation formulas for the LSTM cell

Example: peephole weight output gate



The error gradient at the input of cell state

$$\begin{aligned}\delta_{\chi}^t &= \frac{\partial O}{\partial o_{\chi}^t} = w_{\chi\omega} \delta_{\omega} + \delta_C \\ &\quad + o_{\phi}^{t+1} \delta_{\chi}^{t+1} \\ &\quad + w_{\chi\iota} \delta_{\iota}^{t+1} + w_{\chi\phi} \delta_{\phi}^{t+1}\end{aligned}$$

Back-propagation

- ▶ All error gradients can be constructed in a similar way
- ▶ The result is a memory cell that can be trained via back-propagation