

Why do we like LSTMs?

Tiago Pimentel



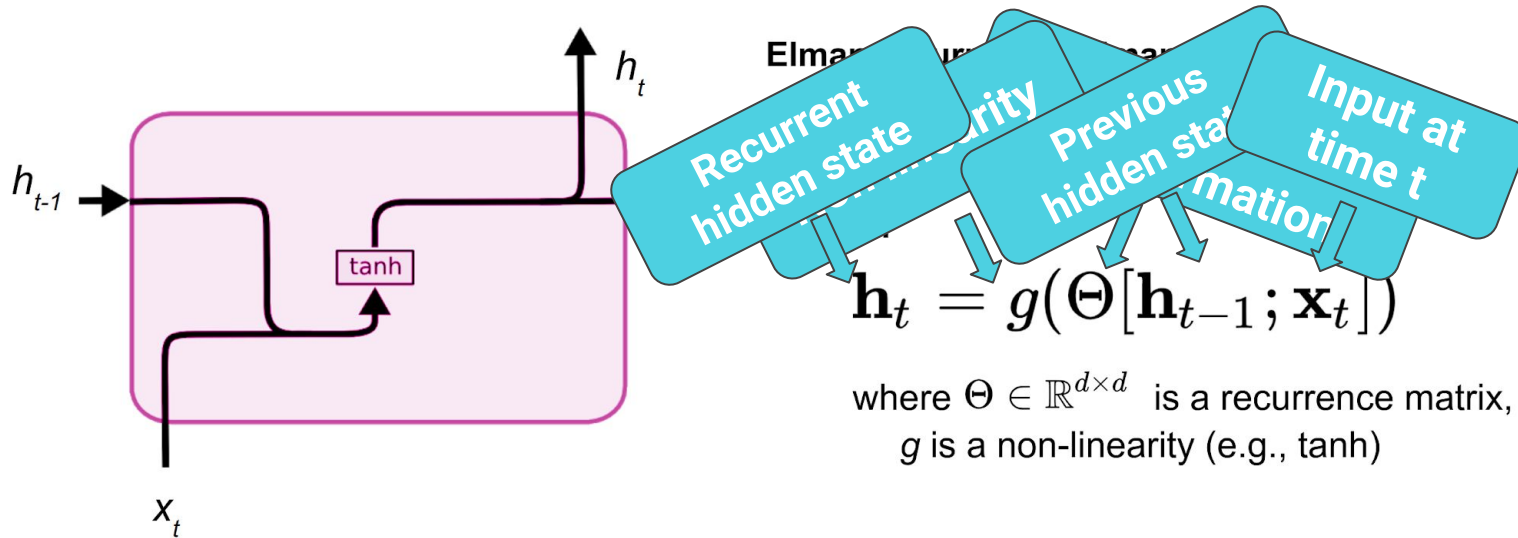
UNIVERSITY OF
CAMBRIDGE

ETH zürich

3 Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



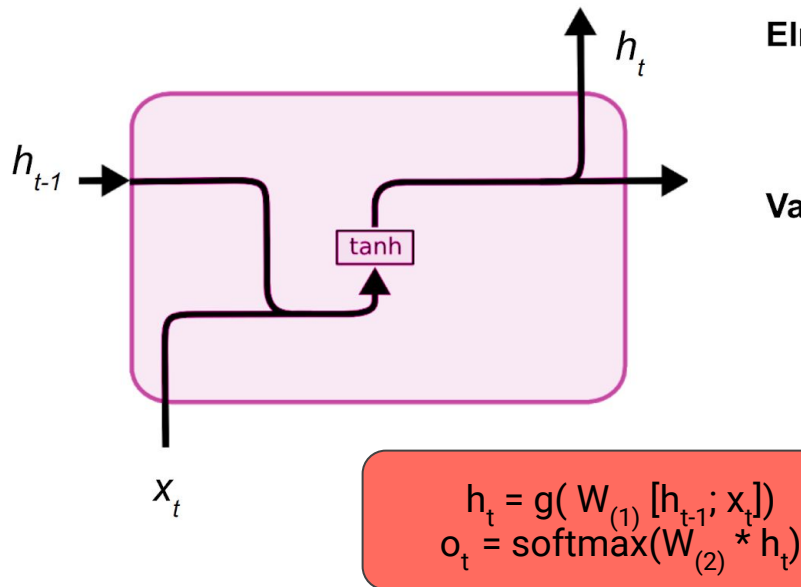
from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



3 Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



Elman recurrence (Elman, 1990):

$$\mathbf{h}_t = g(\Theta \mathbf{h}_{t-1} + \mathbf{x}_t)$$

Variant:

$$\mathbf{h}_t = g(\Theta[\mathbf{h}_{t-1}; \mathbf{x}_t])$$

where $\Theta \in \mathbb{R}^{d \times d}$ is a recurrence matrix,
 g is a non-linearity (e.g., tanh)

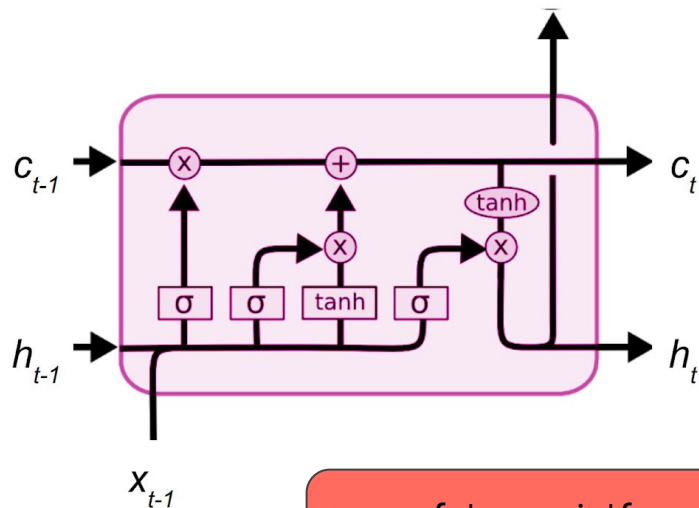


from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



$$c_t = f_t * c_{t-1} + i_t * \text{func}_1(c_{t-1}, x_t)$$

$$o_t = \text{softmax}(W_{(2)} * \text{func}_2(c_t, x_t))$$

Previous hidden state

Modified input at time t

forget gate

input gate

update candidate

update cell state

output gate

update hidden state

$$f_t = \sigma(\Theta_f[h_{t-1}; x_t] + b_f)$$

$$i_t = \sigma(\Theta_i[h_{t-1}; x_t] + b_i)$$

$$\tilde{c}_t = \tanh(\Theta_c[h_{t-1}; x_t] + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

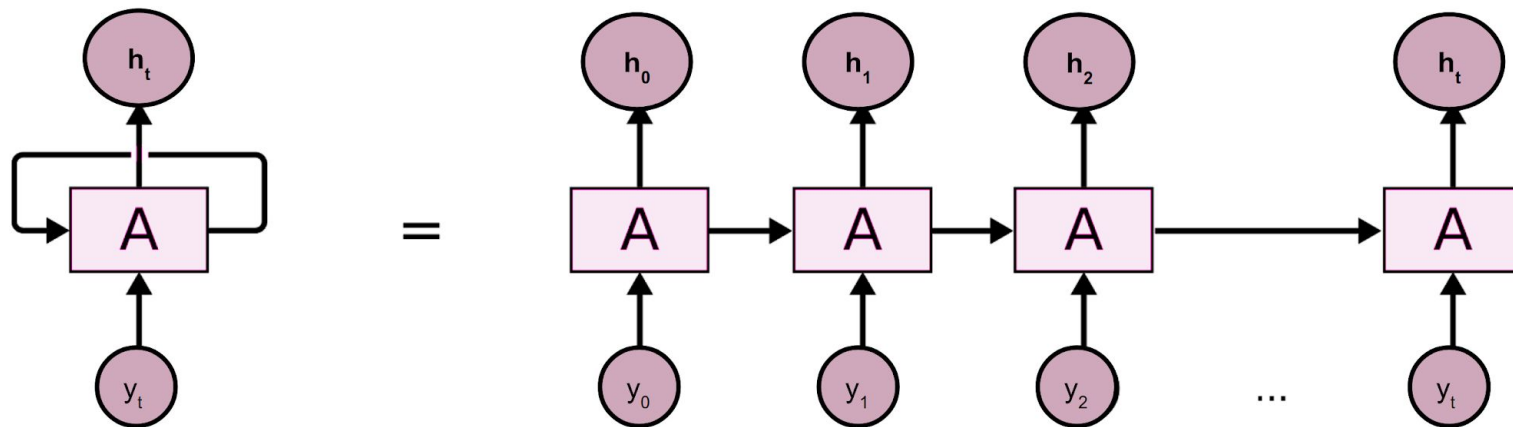
$$o_t = \sigma(\Theta_o[h_{t-1}; x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$



from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3 Backpropagation Through Time is Just Backpropagation



- Here A is whatever RNN cell we want to use (e.g., LSTM, RNN, etc.). Each timestep yields (i) an output and (ii) a recurrent connection.
- Backpropagating RNNs is the same as backpropagating through any neural network, except the parameters are shared across timesteps.
- This same idea can be used to, e.g., tie embeddings or reuse a filter in a CNN

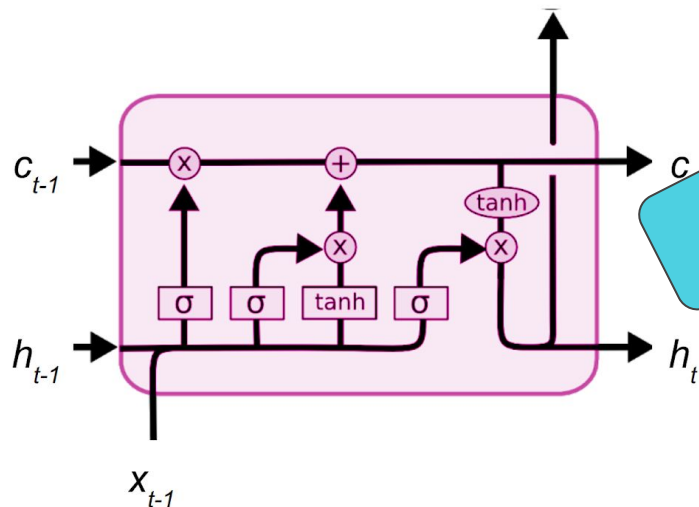
from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3 Recurrent Neural Networks

Original

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



Introduced later in 1999

$$\begin{aligned}
 \mathbf{f}_t &= \sigma(\Theta_f[\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_f) && \text{forget gate} \\
 \mathbf{i}_t &= \sigma(\Theta_i[\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_i) && \text{input gate} \\
 \tilde{\mathbf{c}}_t &= \tanh(\Theta_c[\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_c) && \text{update candidate} \\
 \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t && \text{update cell state} \\
 \mathbf{o}_t &= \sigma(\Theta_o[\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_o) && \text{output gate} \\
 \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t) && \text{update hidden state}
 \end{aligned}$$

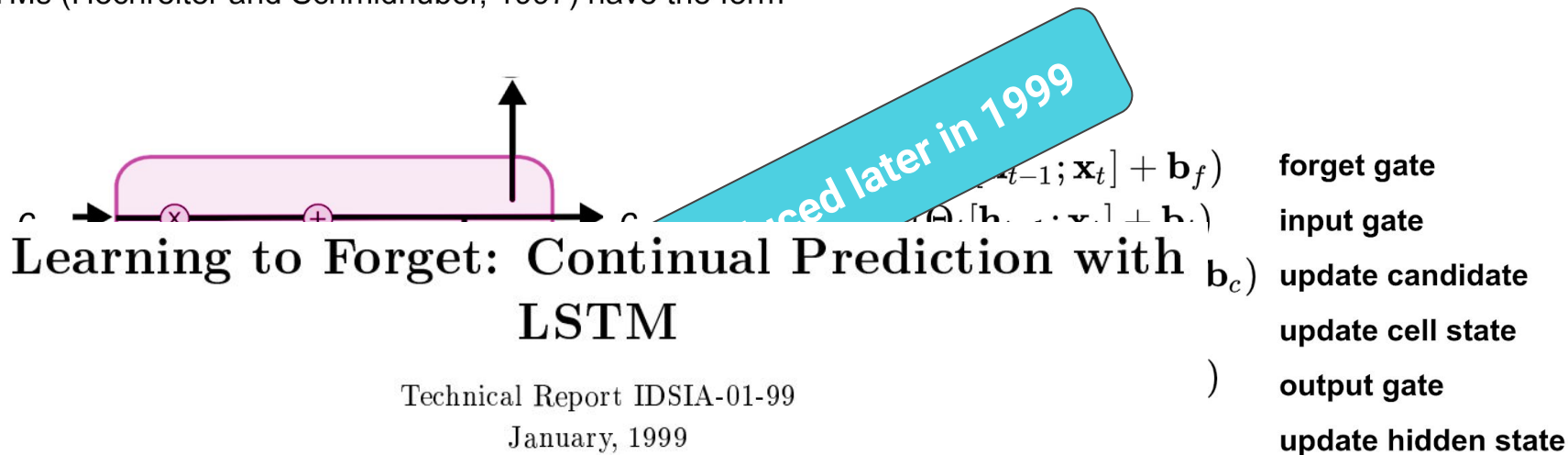


from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



Felix A. Gers
felix@idsia.ch

Jürgen Schmidhuber
juergen@idsia.ch

Fred Cummins
fred@idsia.ch

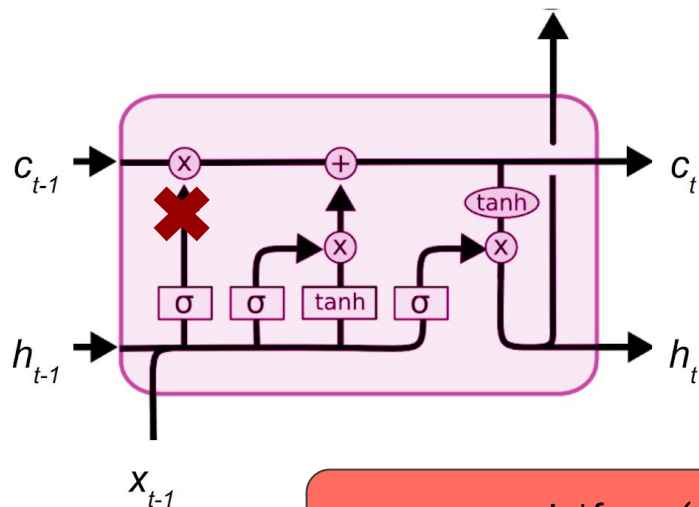
IDSIA, Corso Elvezia 36
6900 Lugano, Switzerland
www.idsia.ch



3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



$$\begin{aligned} c_t &= c_{t-1} + i_t * \text{func}_1(c_{t-1}, x_t) \\ o_t &= \text{softmax}(W_{(2)} * \text{func}_2(c_t, x_t)) \end{aligned}$$

~~$$f_t = \sigma(\Theta_f[h_{t-1}; x_t] + b_f)$$~~

forget gate

$$i_t = \sigma(\Theta_i[h_{t-1}; x_t] + b_i)$$

input gate

$$\tilde{c}_t = \tanh(\Theta_c[h_{t-1}; x_t] + b_c)$$

update candidate

~~$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$~~

update cell state

$$o_t = \sigma(\Theta_o[h_{t-1}, x_t] + b_o)$$

output gate

$$h_t = o_t \odot \tanh(c_t)$$

update hidden state

Neural Network
Layer

Pointwise
Operation

Vector
Transfer

Concatenate

Copy

from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Let's go back to the Elman RNN

Which functions can Elman's RNNs represent?

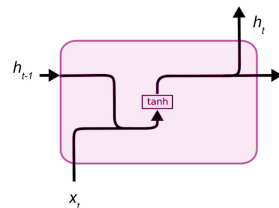
In theory, everything. They are turing complete!

-> "On the computational power of neural nets", Siegelmann & Sontag (1995)

③ Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



Elman recurrence (Elman, 1990):

$$\mathbf{h}_t = g(\Theta \mathbf{h}_{t-1} + \mathbf{x}_t)$$

Variant:

$$\mathbf{h}_t = g(\Theta[\mathbf{h}_{t-1}; \mathbf{x}_t])$$

where $\Theta \in \mathbb{R}^{d \times d}$ is a recurrence matrix,
 g is a non-linearity (e.g., tanh)

from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

ETH zürich



Let's go back to the Elman RNN

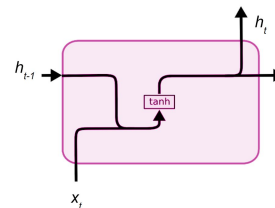
So why would we need something else?

Because some functions might not be **learnable** in practice.

③ Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



Elman recurrence (Elman, 1990):

$$\mathbf{h}_t = g(\Theta \mathbf{h}_{t-1} + \mathbf{x}_t)$$

Variant:

$$\mathbf{h}_t = g(\Theta[\mathbf{h}_{t-1}; \mathbf{x}_t])$$

where $\Theta \in \mathbb{R}^{d \times d}$ is a recurrence matrix,
 g is a non-linearity (e.g., tanh)

from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

ETH zürich



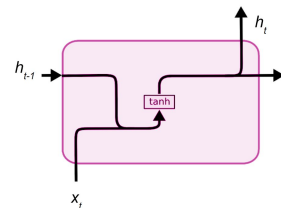
Let's go back to the Elman RNN

Let's unfold this RNN on time

3 Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



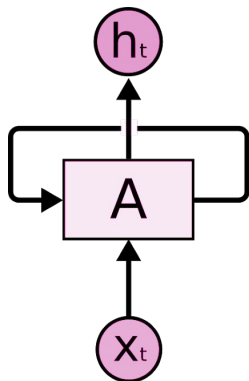
Elman recurrence (Elman, 1990):

$$\mathbf{h}_t = g(\Theta \mathbf{h}_{t-1} + \mathbf{x}_t)$$

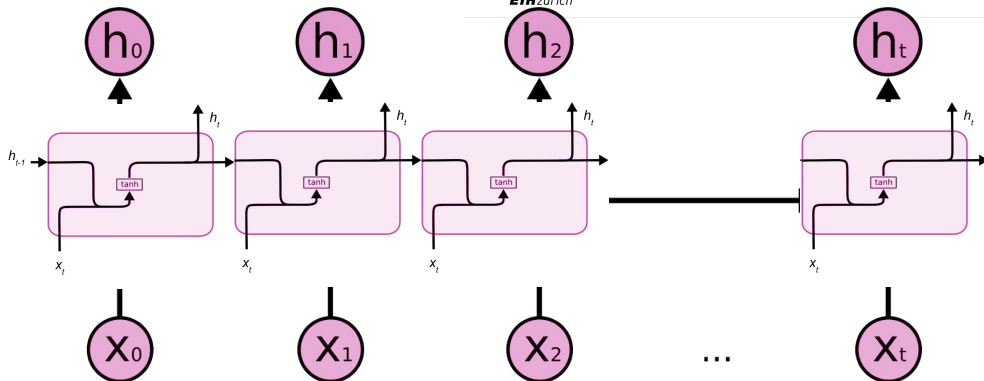
Variant:

$$\mathbf{h}_t = g(\Theta[\mathbf{h}_{t-1}; \mathbf{x}_t])$$

where $\Theta \in \mathbb{R}^{d \times d}$ is a recurrence matrix,
 g is a non-linearity (e.g., tanh)



=



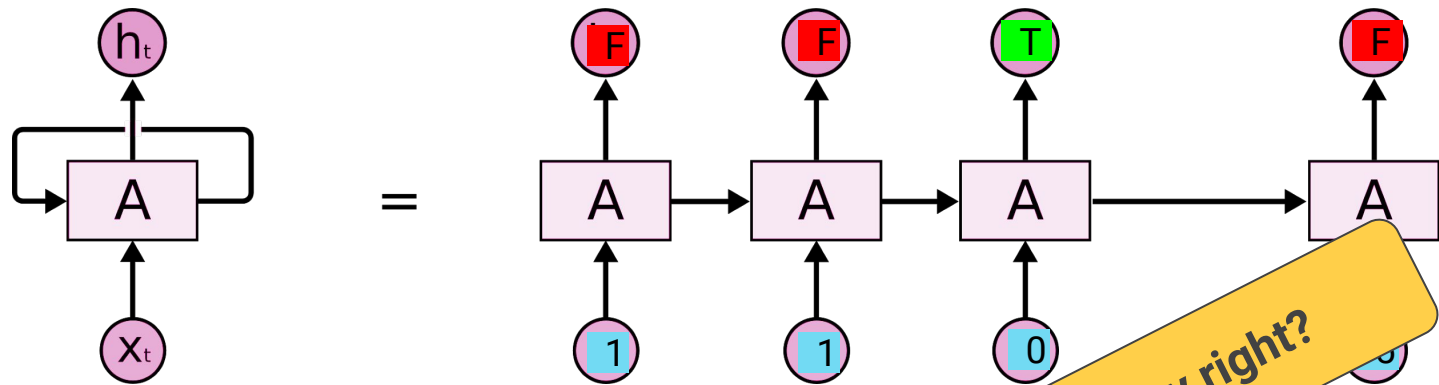
from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

ETH zürich



A simple example

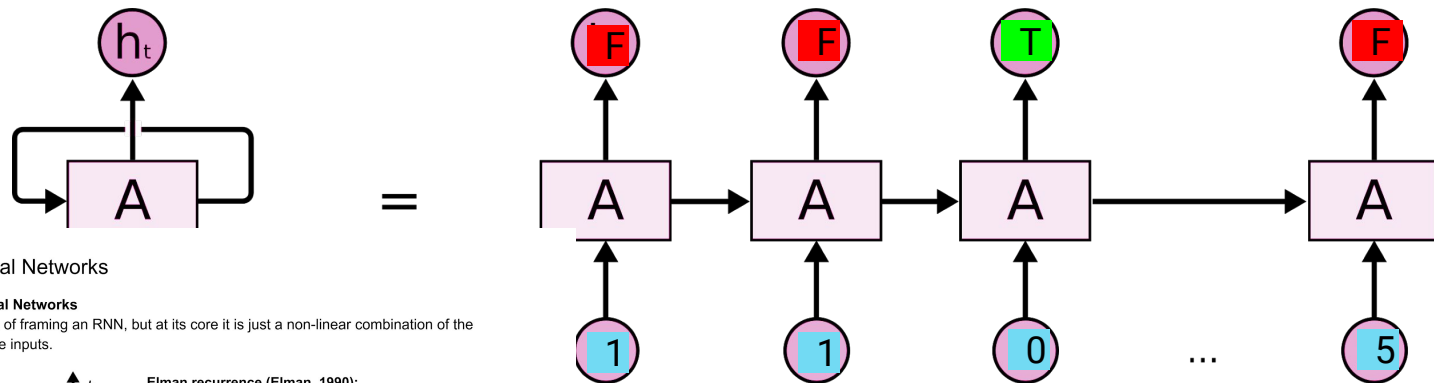
Our task is to identify the number 0 in the input



Pretty easy right?

A simple example

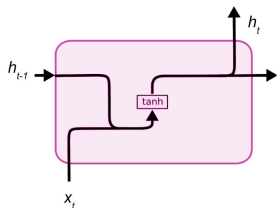
Our task is to identify the number 0 in the input



3 Recurrent Neural Networks

(Vanilla) Recurrent Neural Networks

- There are many ways of framing an RNN, but at its core it is just a non-linear combination of the recurrent state and the inputs.



Elman recurrence (Elman, 1990):

$$\mathbf{h}_t = g(\Theta \mathbf{h}_{t-1} + \mathbf{x}_t)$$

Variant:

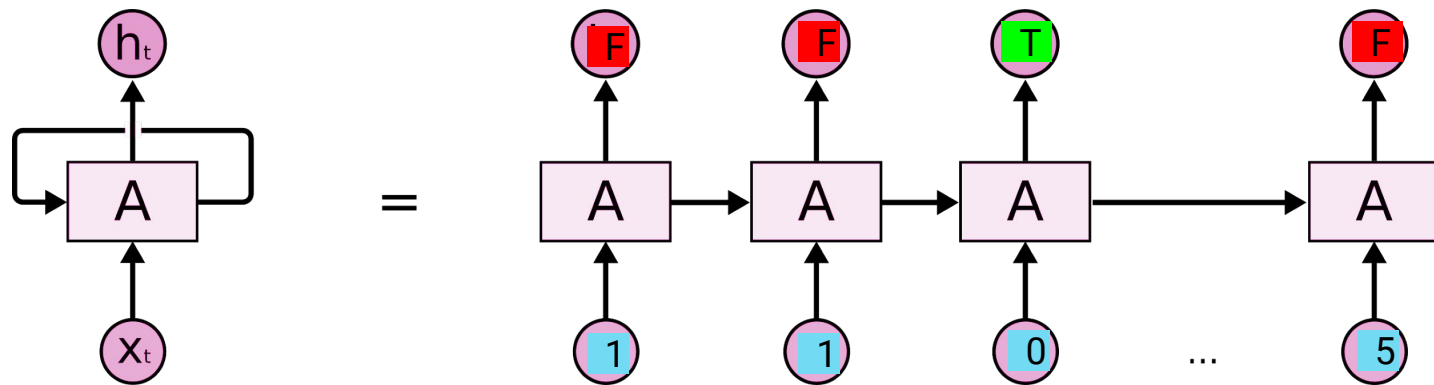
$$\mathbf{h}_t = g(\Theta [\mathbf{h}_{t-1}; \mathbf{x}_t])$$

where $\Theta \in \mathbb{R}^{d \times d}$ is a recurrence matrix,
 g is a non-linearity (e.g., tanh)

$$\begin{aligned} \mathbf{h}_t &= g(W_{(1)} [\mathbf{h}_{t-1}; \mathbf{x}_t]) \\ \mathbf{o}_t &= \text{softmax}(W_{(2)}^* \mathbf{h}_t) \end{aligned}$$

A simple example

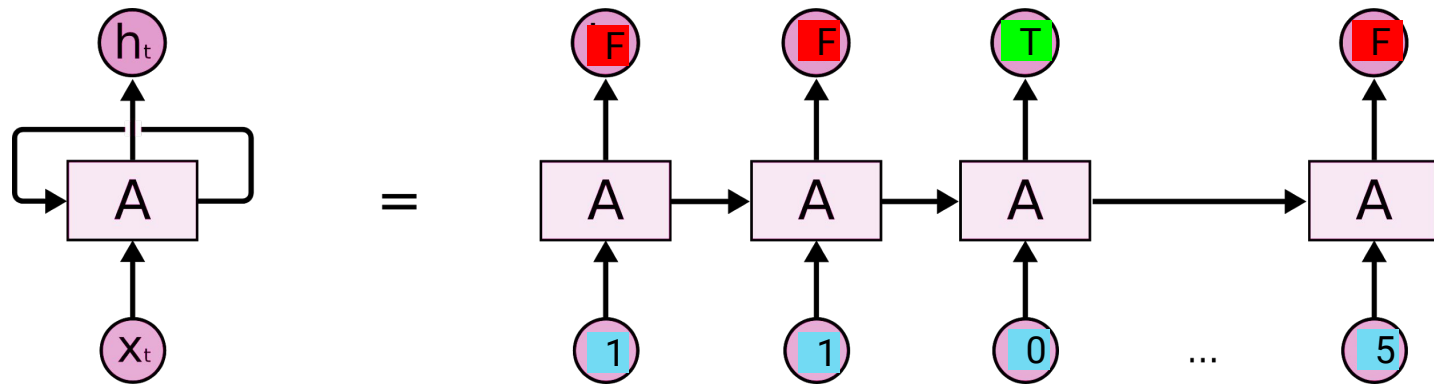
Our task is to identify the number 0 in the input



$$h_t = g(W_{(1)} [h_{t-1}; x_t])$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$

A simple example

Our task is to identify the number 0 in the input

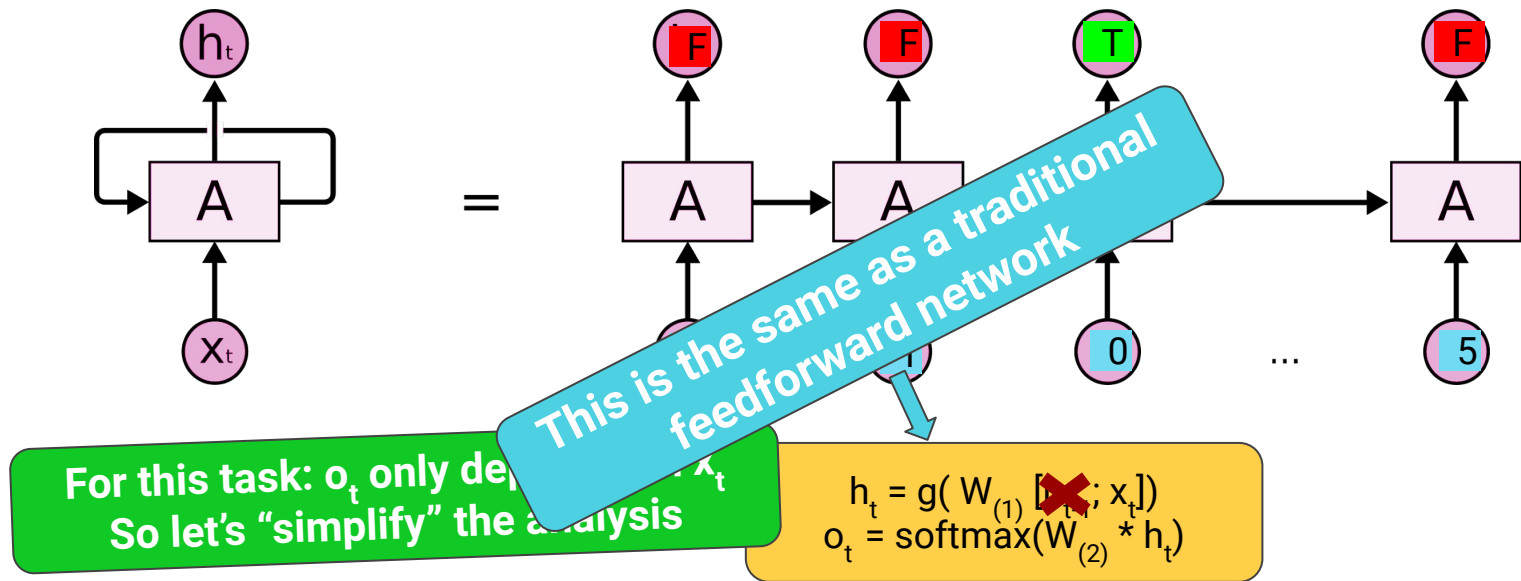


For this task: o_t only depends on x_t

$$h_t = g(W_{(1)} [h_{t-1}; x_t])$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$

A simple example

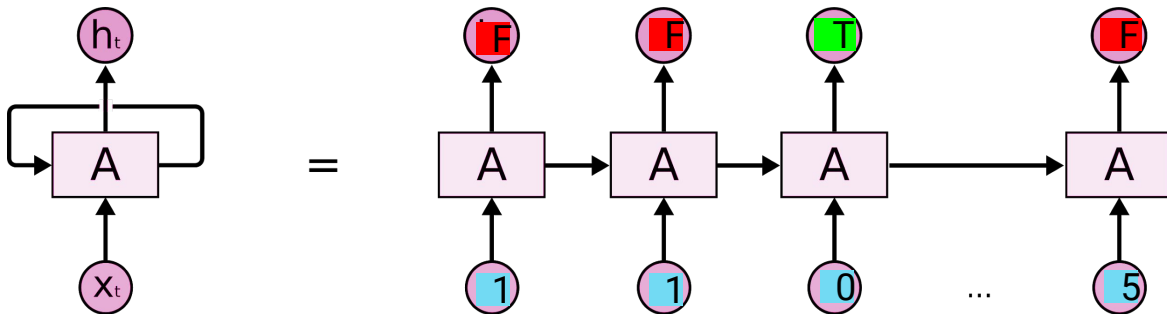
Our task is to identify the number 0 in the input



A simple example

Let's estimate the gradients

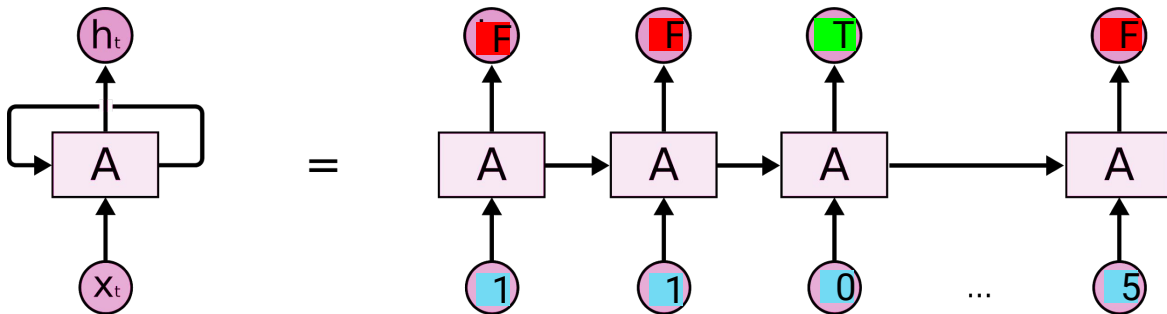
$$h_t = g(W_{(1)} [x_t]; x_t)$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$



A simple example

Let's estimate the gradients

$$h_t = g(W_{(1)} x_t)$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$

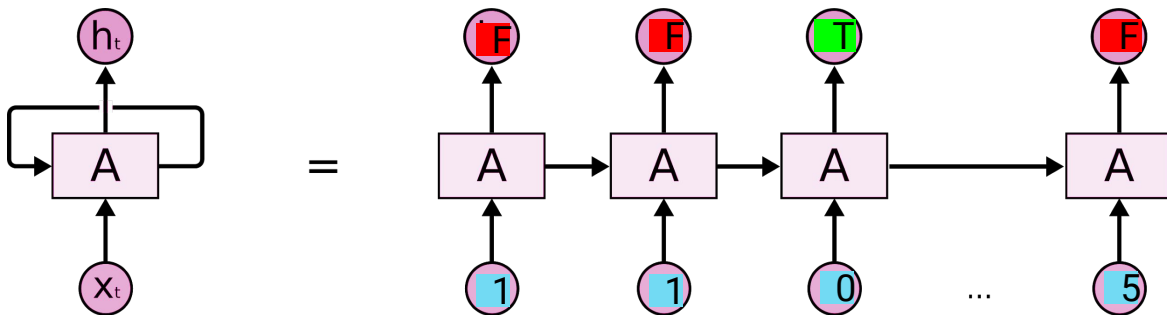


A simple example

Let's estimate the gradients

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_t)$$

$$h_t = g(W_{(1)} x_t)$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$



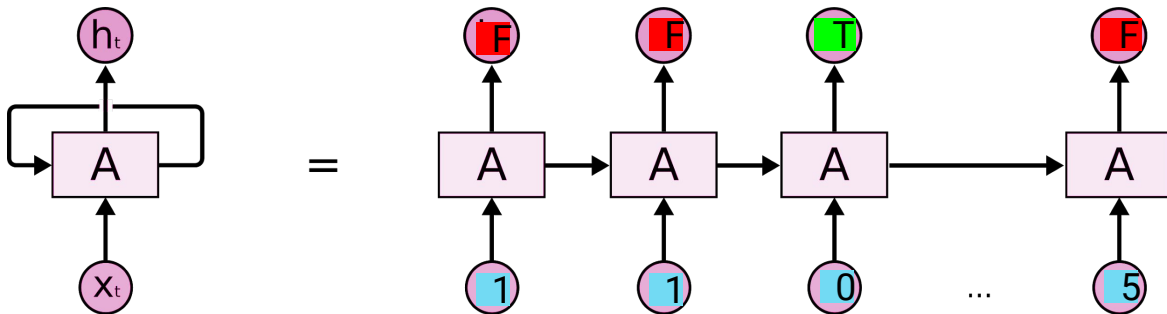
A simple example

Let's estimate the gradients

$$h_t = g(W_{(1)} x_t)$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_t)$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$



A simple example

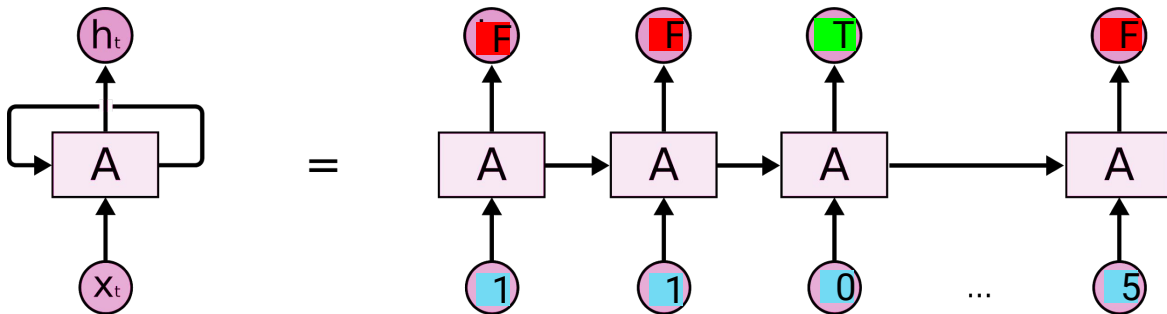
Let's estimate the gradients

$$h_t = g(W_{(1)} x_t)$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_t)$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

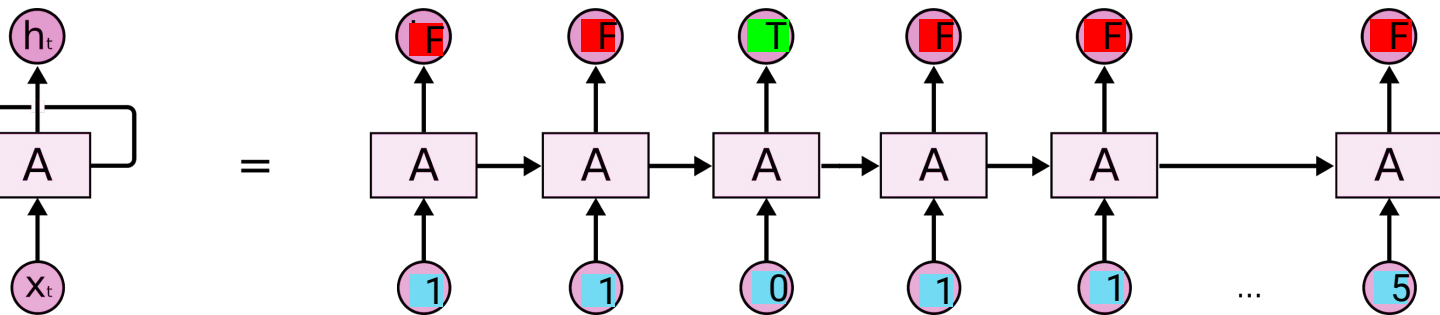
$$dh_t / dx_t = g'(W_{(1)} * x_t) * W_{(1)}$$



A less simple example

Let's change the task a little.

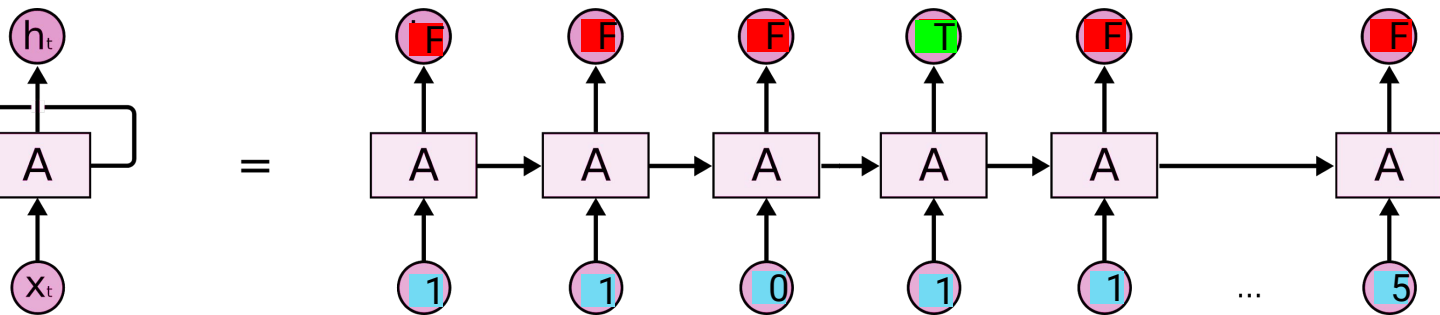
We will add some time lag in the response



A less simple example

Let's change the task a little.

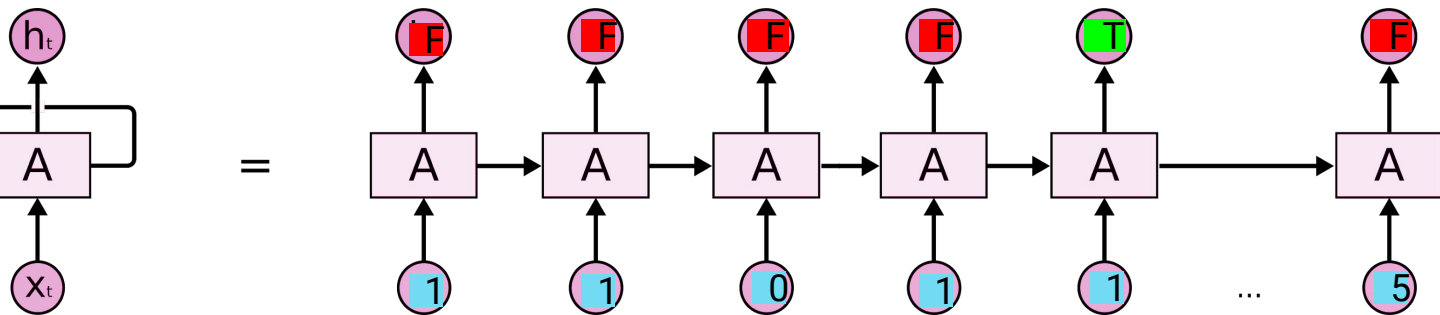
We will add some time lag in the response



A less simple example

Let's change the task a little.

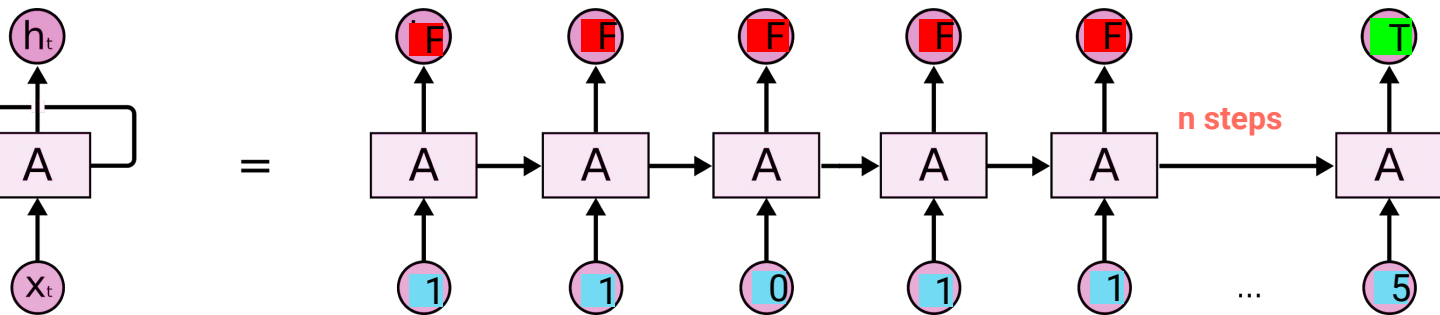
We will add some time lag in the response



A less simple example

Let's change the task a little.

We will add some time lag in the response

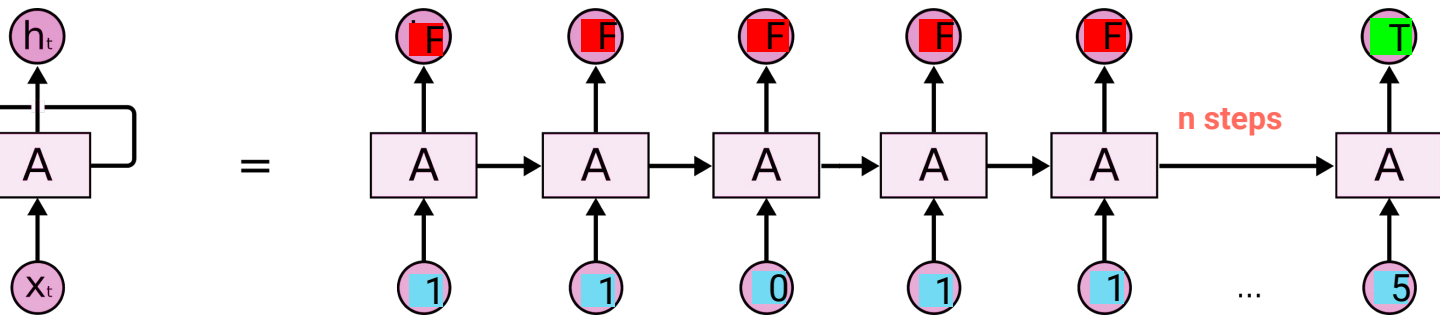


A less simple example

How will the RNN learn this?

We will ignore
this term now

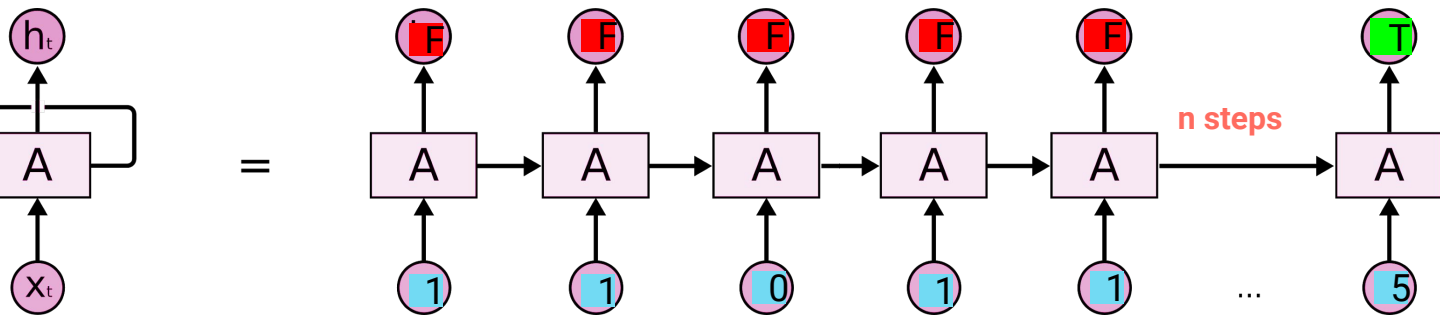
$$h_t = g(W_{(1)} [h_{t-1} \times])$$
$$o_t = \text{softmax}(W_{(2)} * h_t)$$



A less simple example

How will the RNN learn this?

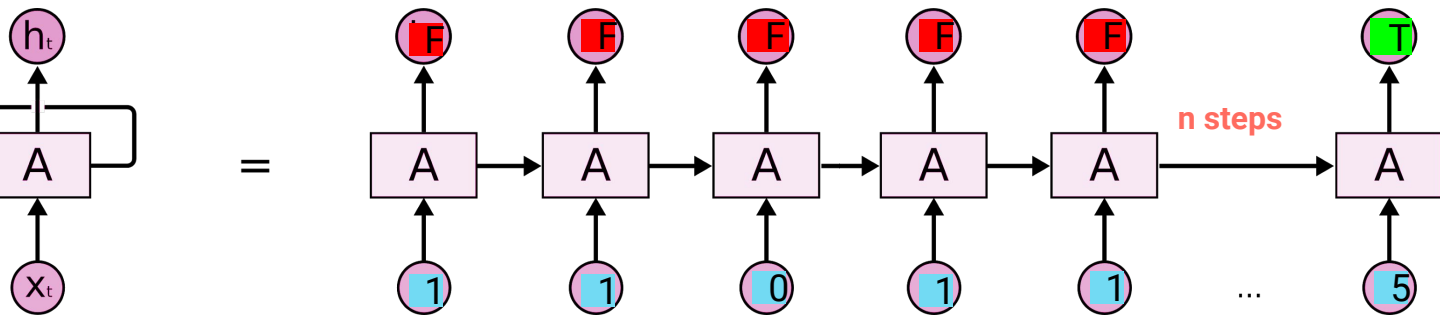
$$o_t = \text{softmax}(W_{(2)} * h_t)$$
$$h_t = g(W_{(1)} [h_{t-1}; x_t])$$



A less simple example

How will the RNN learn this?

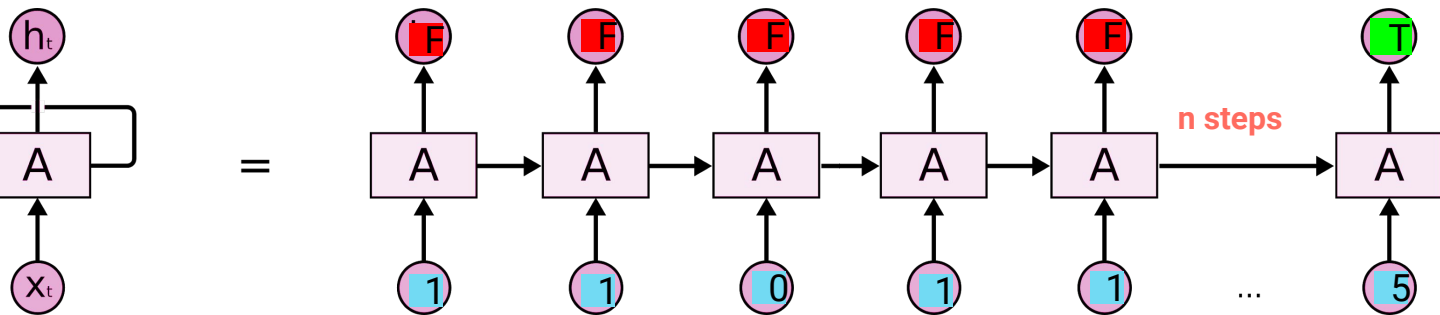
$$o_t = \text{softmax}(W_{(2)} * h_t)$$
$$h_t = g(W_{(1)} [h_{t-1} \text{ } x_t])$$



A less simple example

How will the RNN learn this?

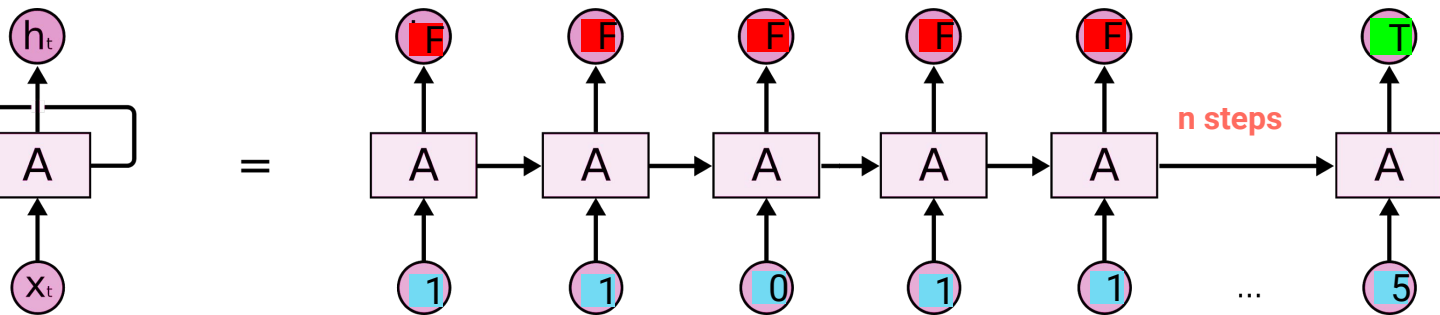
$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \end{aligned}$$



A less simple example

How will the RNN learn this?

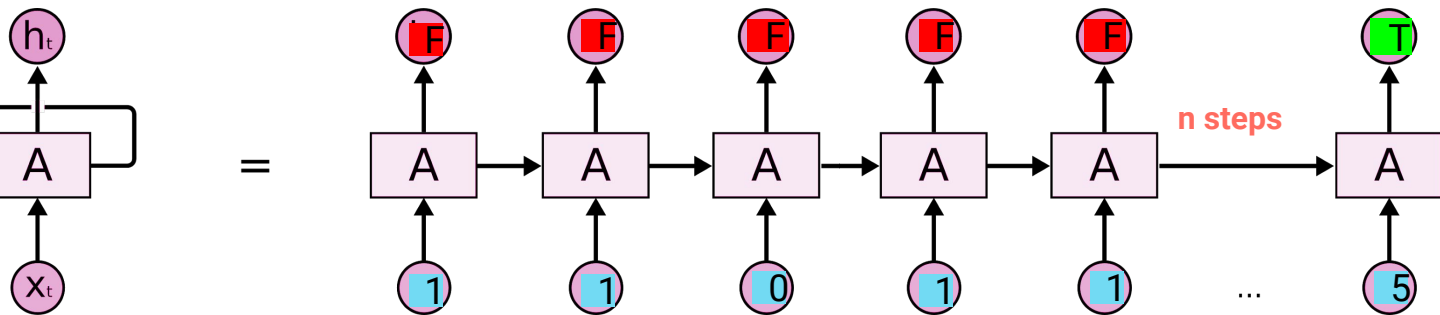
$$\begin{aligned}o_t &= \text{softmax}(W_{(2)} * h_t) \\h_t &= g(W_{(1)} [h_{t-1}]) \\h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\h_{t-2} &= g(W_{(1)} [h_{t-3}])\end{aligned}$$



A less simple example

How will the RNN learn this?

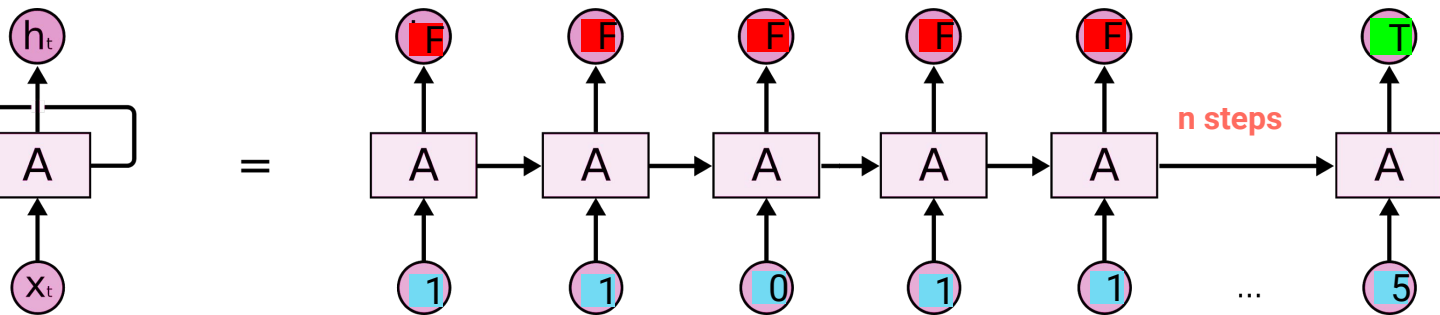
$$\begin{aligned}o_t &= \text{softmax}(W_{(2)} * h_t) \\h_t &= g(W_{(1)} [h_{t-1}]) \\h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\&\dots\end{aligned}$$



A less simple example

How will the RNN learn this?

$$\begin{aligned}o_t &= \text{softmax}(W_{(2)} * h_t) \\h_t &= g(W_{(1)} [h_{t-1}]) \\h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\&\dots \\h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t])\end{aligned}$$



A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

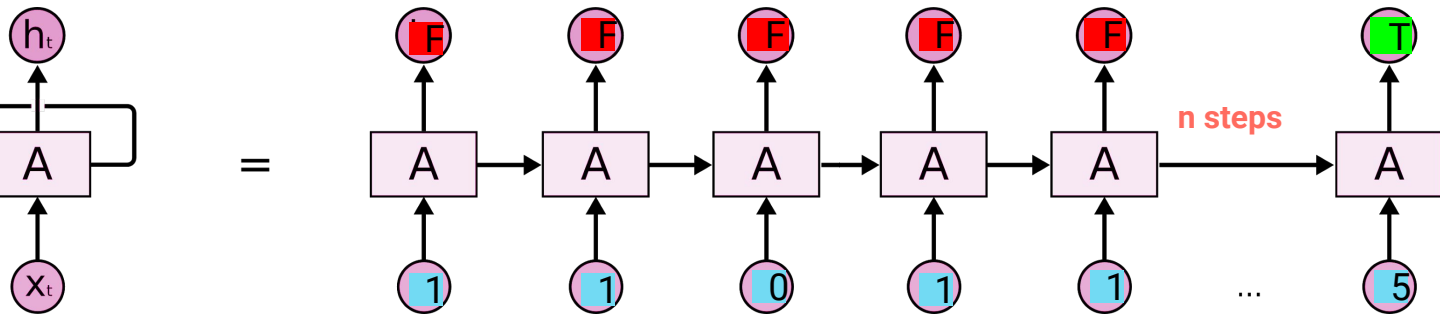
$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

$$\dots$$

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}; x_t])$$



A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$dh_{t-1} / dx_{t-n} = ???$$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$dh_{t-1} / dx_{t-n} = (dh_{t-1} / dh_{t-2}) * (dh_{t-2} / dx_{t-n})$$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\dots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$dh_{t-1} / dx_{t-n} = (dh_{t-1} / dh_{t-2}) * (dh_{t-2} / dx_{t-n})$$

$$dh_{t-2} / dx_{t-n} = (dh_{t-2} / dh_{t-3}) * (dh_{t-3} / dx_{t-n})$$

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

$$\dots$$

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}; x_t])$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$dh_{t-1} / dx_{t-n} = (dh_{t-1} / dh_{t-2}) * (dh_{t-2} / dx_{t-n})$$

$$dh_{t-2} / dx_{t-n} = (dh_{t-2} / dh_{t-3}) * (dh_{t-3} / dx_{t-n})$$

...

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

$$\dots$$

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}; x_t])$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

...

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}; x_t])$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dh_{t-1}) * (dh_{t-1} / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dh_{t-1} = g'(W_{(1)} * h_{t-1}) * W_{(1)}$$

$$dh_{t-1} / dx_{t-n} = (dh_{t-1} / dh_{t-2}) * (dh_{t-2} / dx_{t-n})$$

$$dh_{t-2} / dx_{t-n} = (dh_{t-2} / dh_{t-3}) * (dh_{t-3} / dx_{t-n})$$

...

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

$$\dots$$

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}; x_t])$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

Let's estimate the gradients again!

$$do_t / dx_t = (do_t / dh_t) * (dh_t / dx_{t-n})$$

$$do_t / dh_t = \text{softmax}'(W_{(2)} * h_t) * W_{(2)}$$

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

$$o_t = \text{softmax}(W_{(2)} * h_t)$$

$$h_t = g(W_{(1)} [h_{t-1}])$$

$$h_{t-1} = g(W_{(1)} [h_{t-2}])$$

$$h_{t-2} = g(W_{(1)} [h_{t-3}])$$

...

$$h_{t-n+1} = g(W_{(1)} [h_{t-n}]; x_t)$$

What will happen if **n** is large?

A less simple example

dh_t / dh_{t-1}

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if n is large?

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if **n** is large? It depends

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if n is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if n is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if **n** is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| < 1$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if **n** is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| < 1$
 - Gradients will exponentially shrink and “vanish”

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if **n** is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| < 1$
 - Gradients will exponentially shrink and “vanish”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| = 1$

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if **n** is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| < 1$
 - Gradients will exponentially shrink and “vanish”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| = 1$
 - Gradients can flow through “infinite” time steps

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1}]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2}]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3}]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n}; x_t]) \end{aligned}$$

A less simple example

$$dh_t / dh_{t-1}$$

$$dh_t / dx_{t-n} = \prod_{i=0}^n [g'(W_{(1)} * h_{t-i}) * W_{(1)}] * (dh_{t-n} / dx_{t-n})$$

What will happen if n is large? It depends

- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| > 1$
 - Gradients will exponentially grow and “explode”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| < 1$
 - Gradients will exponentially shrink and “vanish”
- $|g'(W_{(1)} * h_{t-i}) * W_{(1)}| = 1$
 - Gradients can flow through “infinite” time steps

$$\begin{aligned} o_t &= \text{softmax}(W_{(2)} * h_t) \\ h_t &= g(W_{(1)} [h_{t-1} \text{ }]) \\ h_{t-1} &= g(W_{(1)} [h_{t-2} \text{ }]) \\ h_{t-2} &= g(W_{(1)} [h_{t-3} \text{ }]) \\ &\vdots \\ h_{t-n+1} &= g(W_{(1)} [h_{t-n} \text{ } x_t]) \end{aligned}$$

So why LSTMs?

Why LSTMs?

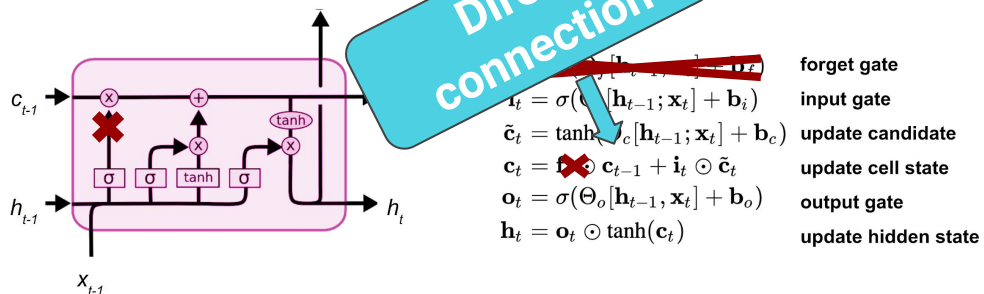
We want gradients to flow!

$$\text{i.e. } dh_t / dh_{t-1} = 1$$

3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Why LSTMs?

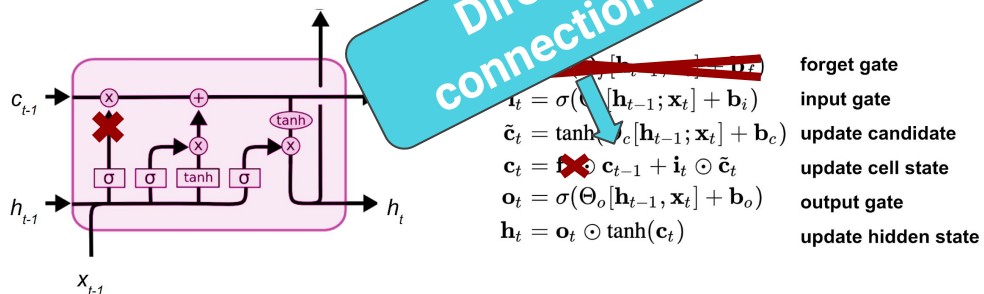
We want gradients to flow!

$$\text{i.e. } dc_t / dc_{t-1} = 1$$

3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form



from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Why LSTMs?

We want gradients to f

$$\text{i.e. } dc_t / dc_{t-1} = 1$$

Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun
Microsoft Research
{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

Same trick as ResNet, but ~20 years before

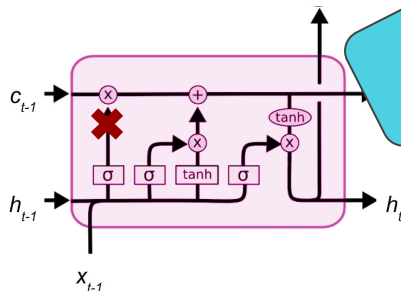
Does this remind you of something else?

3 Recurrent Neural Networks

Long short-term memory (LSTM)

LSTMs (Hochreiter and Schmidhuber, 1997) have the form

Direct connection



forget gate
input gate
update candidate
update cell state
output gate
update hidden state

$$i_t = \sigma(\Theta_i[h_{t-1}; x_t] + b_i)$$

$$\tilde{c}_t = \tanh(\Theta_c[h_{t-1}; x_t] + b_c)$$

$$c_t = c_{t-1} \odot i_t + \tilde{c}_t \odot i_t$$

$$o_t = \sigma(\Theta_o[h_{t-1}; x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$

$$c_t = c_{t-1} + i_t * \text{func}_1(c_{t-1}, x_t)$$

$$o_t = \text{softmax}(W_{(2)} * \text{func}_2(c_t, x_t))$$

from Christopher Olah's blog, <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>