# Forget to remember
# Remember to forget

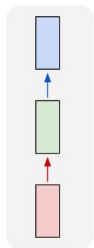# Long Short Term Memories and Gated Recurrent Units

## Kate Farrahi

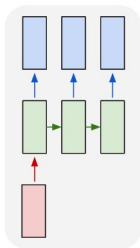Vision, Learning and Control
University of Southampton

Some of the slides, images and animations used here were originally designed by Jonathon Hare and Adam Prügel-Bennett.
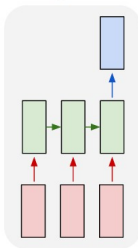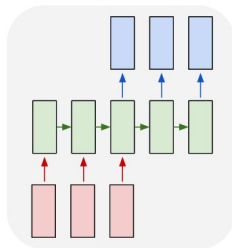
# Recap: RNN
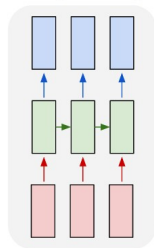


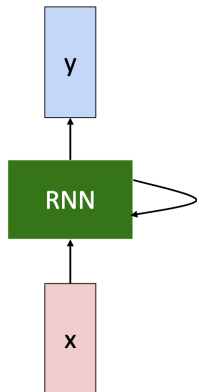| one to one | one to many | many to one | many to many | many to many |

Example Applications: image classification, image captioning, video classification, machine translation, language modeling

# Recap: Vanilla or Elman RNN

Excluding bias terms



$$h_t = tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Image from https://web.eecs.umich.edu

RNN Computational Graph (Many to Many)

Image from https://web.eecs.umich.edu

# Example: Character-level language modelling



Image from `http://karpathy.github.io/2015/05/21/rnn-effectiveness/`

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients



$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

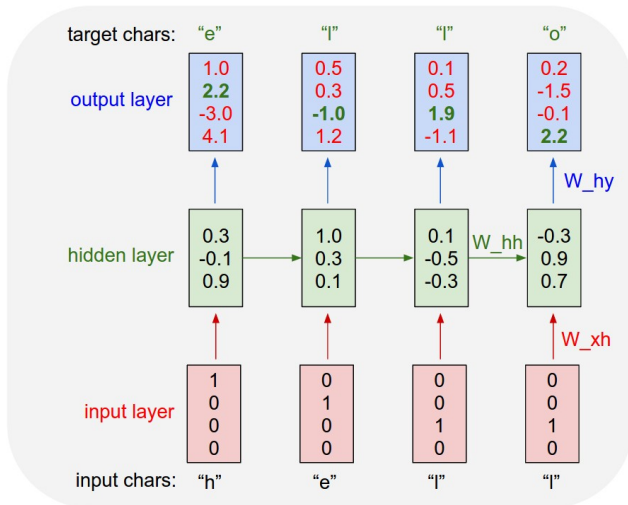# Vanishing and Exploding Gradients

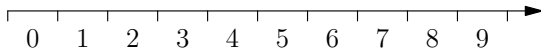$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients

$$y(t) = w_1 \big( x(t) + w_2 \, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients



$$y(t) = w_1 \big( x(t) + w_2 \, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients
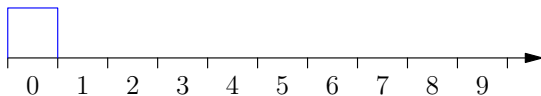
$$y(t) = w_1 \big( x(t) + w_2 \, y(t-1) \big)$$
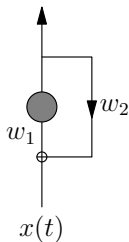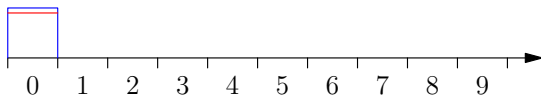
$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients

$$y(t) = w_1 \big(x(t) + w_2\, y(t-1)\big)$$

$w_1 = w_2 = 0.9$

$$y(t) = w_1 \big( x(t) + w_2 \, y(t-1) \big)$$

$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients

$$y(t) = w_1\big(x(t) + w_2\,y(t-1)\big)$$

$$w_1 = w_2 = 0.9$$
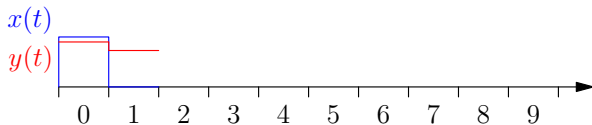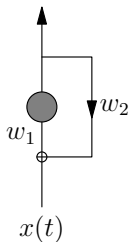
$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$
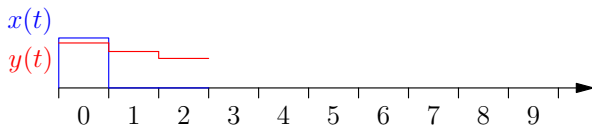
$$w_1 = w_2 = 0.9$$

# Vanishing and Exploding Gradients

$$y(t) = w_1\big(x(t) + w_2\, y(t-1)\big)$$

$w_1 = w_2 = 0.9$

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$w_1 = w_2 = 1.1$$

$$y(t) = w_1 \big(x(t) + w_2\, y(t-1)\big)$$



$$w_1 = w_2 = 1.1$$

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$$w_1 = w_2 = 1.1$$

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$
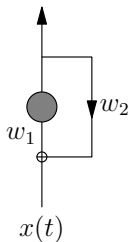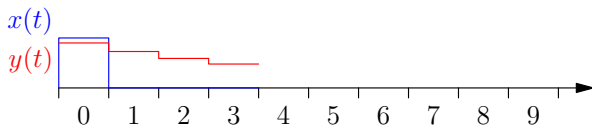
$$w_1 = w_2 = 1.1$$

# Vanishing and Exploding Gradients

$$y(t) = w_1 \big( x(t) + w_2 \, y(t-1) \big)$$

$$w_1 = w_2 = 1.1$$

# Vanishing and Exploding Gradients



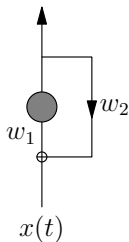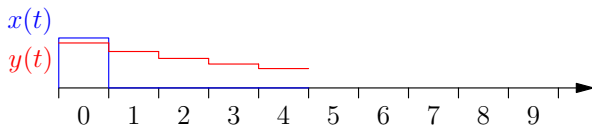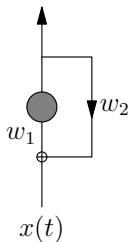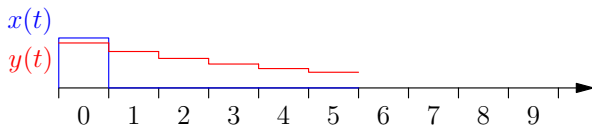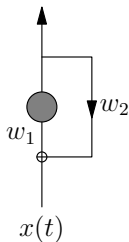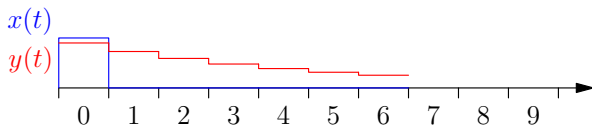$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$w_1 = w_2 = 1.1$

$y(t)$
$x(t)$

$$y(t) = w_1 \big( x(t) + w_2\, y(t-1) \big)$$

$w_1 = w_2 = 1.1$

$y(t)$
$x(t)$

$w_2$
$w_1$
$x(t)$

$$y(t) = w_1\big(x(t) + w_2\, y(t-1)\big)$$

$w_1 = w_2 = 1.1$

$y(t)$
$x(t)$

# Vanishing and Exploding Gradients



Image from https://web.eecs.umich.edu

# LSTM Architecture

- LSTMs (long-short term memory) were designed to solve this problem

# LSTM Architecture

- LSTMs (long-short term memory) were designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\boldsymbol{c}(t) = \boldsymbol{c}(t-1)$$

# LSTM Architecture

- LSTMs (long-short term memory) were designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\boldsymbol{c}(t) = \boldsymbol{c}(t-1)$$

- Sometimes we have to forget and sometimes we have to change a memory

# LSTM Architecture

- LSTMs (long-short term memory) were designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\boldsymbol{c}(t) = \boldsymbol{c}(t-1)$$

- Sometimes we have to forget and sometimes we have to change a memory
- To do this we should use 'gates' that saturate at 0 and 1

# LSTM Architecture

- LSTMs (long-short term memory) were designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\boldsymbol{c}(t) = \boldsymbol{c}(t-1)$$

- Sometimes we have to forget and sometimes we have to change a memory
- To do this we should use 'gates' that saturate at 0 and 1
- Sigmoid functions naturally saturate at 0 and 1

# LSTM Architecture

LSTM

# LSTM Architecture

$$\tanh(c_j(t)) \in \{-1, 1\}$$

LSTM

# LSTM Architecture

# LSTM Architecture

$$\boldsymbol{f}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_f \boldsymbol{z}(t) + \boldsymbol{b}_f)$$

$$\boldsymbol{z}(t) = (\boldsymbol{x}(t), \boldsymbol{h}(t-1))$$

$$f_j(t)\, c_j(t-1)$$

$$\sigma_j(\boldsymbol{W}_f \boldsymbol{z}(t) + \boldsymbol{b}_f)\, c_j(t-1)$$

$$\boldsymbol{f}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_f \boldsymbol{z}(t) + \boldsymbol{b}_f)$$

$$\boldsymbol{z}(t) = (\boldsymbol{x}(t), \boldsymbol{h}(t-1))$$

# LSTM Architecture

# LSTM Architecture

# LSTM Architecture

# LSTM Architecture

# LSTM Architecture

# LSTM Architecture

# Update Equations

Initially, for $t = 0$, $\boldsymbol{h}(0) = \boldsymbol{0}$

- Inputs $\boldsymbol{z}(t) = \big(\boldsymbol{x}(t), \, \boldsymbol{h}(t-1)\big)$

## Update Equations

Initially, for $t = 0$, $\boldsymbol{h}(0) = \boldsymbol{0}$

- Inputs $\boldsymbol{z}(t) = \big(\boldsymbol{x}(t),\, \boldsymbol{h}(t-1)\big)$
- Network updates ($\boldsymbol{W}_*$ and $\boldsymbol{b}_*$ are the learnable parameters)

$$\boldsymbol{f}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_f\,\boldsymbol{z}(t) + \boldsymbol{b}_f) \qquad\qquad \boldsymbol{i}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_i\,\boldsymbol{z}(t) + \boldsymbol{b}_i)$$
$$\boldsymbol{g}(t) = \tanh(\boldsymbol{W}_g\,\boldsymbol{z}(t) + \boldsymbol{b}_g) \qquad \boldsymbol{o}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_o\,\boldsymbol{z}(t) + \boldsymbol{b}_o)$$

# Update Equations

Initially, for $t = 0$, $\boldsymbol{h}(0) = \boldsymbol{0}$

- Inputs $\boldsymbol{z}(t) = \big(\boldsymbol{x}(t), \boldsymbol{h}(t-1)\big)$
- Network updates ($\boldsymbol{W}_*$ and $\boldsymbol{b}_*$ are the learnable parameters)

$$\boldsymbol{f}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_f\,\boldsymbol{z}(t) + \boldsymbol{b}_f) \qquad \boldsymbol{i}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_i\,\boldsymbol{z}(t) + \boldsymbol{b}_i)$$
$$\boldsymbol{g}(t) = \tanh(\boldsymbol{W}_g\,\boldsymbol{z}(t) + \boldsymbol{b}_g) \qquad \boldsymbol{o}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_o\,\boldsymbol{z}(t) + \boldsymbol{b}_o)$$

- Long-term memory update

$$\boldsymbol{c}(t) = \boldsymbol{f}(t) \odot \boldsymbol{c}(t-1) + \boldsymbol{g}(t) \odot \boldsymbol{i}(t)$$

# Update Equations

Initially, for $t = 0$, $\boldsymbol{h}(0) = \boldsymbol{0}$

- Inputs $\boldsymbol{z}(t) = \big(\boldsymbol{x}(t), \boldsymbol{h}(t-1)\big)$
- Network updates ($\boldsymbol{W}_*$ and $\boldsymbol{b}_*$ are the learnable parameters)

$$\boldsymbol{f}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_f\,\boldsymbol{z}(t) + \boldsymbol{b}_f) \qquad \boldsymbol{i}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_i\,\boldsymbol{z}(t) + \boldsymbol{b}_i)$$
$$\boldsymbol{g}(t) = \tanh(\boldsymbol{W}_g\,\boldsymbol{z}(t) + \boldsymbol{b}_g) \qquad \boldsymbol{o}(t) = \boldsymbol{\sigma}(\boldsymbol{W}_o\,\boldsymbol{z}(t) + \boldsymbol{b}_o)$$

- Long-term memory update

$$\boldsymbol{c}(t) = \boldsymbol{f}(t) \odot \boldsymbol{c}(t-1) + \boldsymbol{g}(t) \odot \boldsymbol{i}(t)$$

- Output $\boldsymbol{h}(t) = \boldsymbol{o}(t) \odot \tanh(\boldsymbol{c}(t))$

- We can train an LSTM by unwrapping it in time.

- We can train an LSTM by unwrapping it in time.
- Note that it involves four dense layers with sigmoidal (or tanh) outputs.

- We can train an LSTM by unwrapping it in time.
- Note that it involves four dense layers with sigmoidal (or tanh) outputs.
- This means that typically it is very slow to train.

# Training LSTMs

- We can train an LSTM by unwrapping it in time.
- Note that it involves four dense layers with sigmoidal (or tanh) outputs.
- This means that typically it is very slow to train.
- There are a few variants of LSTMs, but all are very similar. The most popular is probably the Gated Recurrent Unit (GRU).

## LSTM Success Stories

- LSTMs have been used to win many competitions in speech and handwriting recognition.
- Major technology companies including Google, Apple, and Microsoft are using LSTMs as fundamental components in products.
- Google used LSTM for speech recognition on the smartphone, for Google Translate.
- Apple uses LSTM for the "Quicktype" function on the iPhone and for Siri.
- Amazon uses LSTM for Amazon Alexa.
- In 2017, Facebook performed some 4.5 billion automatic translations every day using long short-term memory networks[1].

---

[1] https://en.wikipedia.org/wiki/Long_short-term_memory

GRU

# Gated Recurrent Unit (GRU)

- $x(t)$: input vector
- $h(t)$: output vector (and 'hidden state')
- $r(t)$: reset gate vector
- $z(t)$: update gate vector
- $n(t)$: new state vector (before update is applied)
- $W$ and $b$: parameter matrices and biases

# Gated Recurrent Unit (GRU)

Initially, for $t = 0$, $\boldsymbol{h}(0) = \boldsymbol{0}$

$$\boldsymbol{z}(t) = \sigma(\boldsymbol{W}_z(\boldsymbol{x}(t), \boldsymbol{h}(t-1)) + \boldsymbol{b}_z)$$
$$\boldsymbol{r}(t) = \sigma(\boldsymbol{W}_r(\boldsymbol{x}(t), \boldsymbol{h}(t-1)) + \boldsymbol{b}_r)$$
$$\boldsymbol{n}(t) = \tanh(\boldsymbol{W}_n(\boldsymbol{x}(t), r(t) \odot h(t-1)) + \boldsymbol{b}_h)$$
$$\boldsymbol{h}(t) = (1 - \boldsymbol{z}(t)) \odot \boldsymbol{h}(t-1) + \boldsymbol{z}(t) \odot \boldsymbol{n}(t)$$

---

Most implementations follow the original paper and swap $(1 - \boldsymbol{z}(t))$ and $(\boldsymbol{z}(t))$ in the $\boldsymbol{h}(t)$ update; this doesn't change the operation of the network, but does change the interpretation of the update gate, as the gate would have to produce a 0 when an update was to occur, and a 1 when no update is to happen (which is somewhat counter-intuitive)!

# GRU or LSTM?

- GRUs have two gates (reset and update) whereas LSTM has three gates (input/output/forget)
- GRU performance on par with LSTM but computationally more efficient (less operations & weights).
- In general, if you have a very large dataset then LSTMs will likely perform slightly better.
- GRUs are a good choice for smaller datasets.