# Lecture 8: Recurrent Neural Networks

Deep Learning @ UvA

# Previous Lecture

o Word and Language Representations

o From n-grams to Neural Networks

o Word2vec

o Skip-gram

# Lecture Overview

- ○ Recurrent Neural Networks (RNN) for sequences

- ○ Backpropagation Through Time

- ○ Vanishing and Exploding Gradients and Remedies

- ○ RNNs using Long Short-Term Memory (LSTM)

- ○ Applications of Recurrent Neural Networks

# Recurrent Neural Networks



Output

NN Cell State

Recurrent connections

Input

# Sequences

o Next data depend on previous data

o Roughly equivalent to predicting what comes next

$$\Pr(x) = \prod_i \Pr(x_i | x_1, \dots, x_{i-1})$$

What

# Sequences

o Next data depend on previous data

o Roughly equivalent to predicting what comes next

$$\Pr(x) = \prod_i \Pr(x_i | x_1, \ldots, x_{i-1})$$



you can be cool

but never a parrot
wearing a hoodie cool

*What about*

# Sequences

o Next data depend on previous data

o Roughly equivalent to predicting what comes next

$$\Pr(x) = \prod_i \Pr(x_i | x_1, \dots, x_{i-1})$$



What about inputs

# Sequences

o Next data depend on previous data

o Roughly equivalent to predicting what comes next

$$\Pr(x) = \prod_i \Pr(x_i | x_1, \ldots, x_{i-1})$$


you can be cool
but never a parrot
wearing a hoodie cool

What about inputs that appear in sequences, such as text? Could a neural network handle such modalities?

# Why sequences?

# Why sequences?

o Considering small chunks $x_i$ → fewer parameters, easier modelling

o Generalizes well to arbitrary lengths

$$RecurrentModel(\text{I think, therefore, I am!})$$
$$\equiv$$
$$RecurrentModel(\text{Everything is repeated, in a circle. History is a master because it teaches us that it doesn't exist. It's the permutations that matter.})$$

o However, often we pick a "frame" $T$ instead of an arbitrary length

$$\Pr(x) = \prod_i \Pr(x_i | x_{i-T}, \dots, x_{i-1})$$

# What a sequence *really* is?

o Data inside a sequence are non i.i.d.
  ◦ Identically, independently distributed

o The next "word" depends on the previous "words"
  ◦ Ideally on all of them

o We need **context,** and we need **memory!**

o How to model context and memory ?

I am Bond , James [ ]

McGuire

Bond

tired

am

!

# What a sequence *really* is?

o Data inside a sequence are non i.i.d.
  ◦ Identically, independently distributed

o The next "word" depends on the previous "words"
  ◦ Ideally on all of them

o We need **context,** and we need **memory!**

o How to model context and memory ?

I am Bond , James [Bond]

McGuire
Bond
tired
am
!

# $x_i \equiv$ One-hot vectors

o A vector with all zeros except for the active dimension

o 12 words in a sequence → 12 One-hot vectors

o After the one-hot vectors apply an embedding (Word2Vec, GloVE)

Vocabulary                    One-hot vectors

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| I | | **1** | | 0 | | 0 | 0 |
| am | am | 0 | am | **1** | am | 0 | am 0 |
| Bond | Bond | 0 | Bond | 0 | Bond | **1** | Bond 0 |
| James | James | 0 | James | 0 | James | 0 | James **1** |
| tired | tired | 0 | tired | 0 | tired | 0 | tired 0 |
| , | , | 0 | , | 0 | , | 0 | , 0 |
| McGuire | McGuire | 0 | McGuire | 0 | McGuire | 0 | McGuire 0 |
| ! | ! | 0 | ! | 0 | ! | 0 | ! 0 |

# Indices instead of one-hot vectors?

o Can't we simply use indices as features?

o No, great solution, because introduces artificial bias between inputs

|  | I | am | James | McGuire |
|---|---|---|---|---|

$$x_{t=1,2,3,4} = \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{matrix}$$

|  | I | am | James | McGuire |
|---|---|---|---|---|
| $q_{t=1,2,3,4} =$ | 1 | 2 | 4 | 7 |

$$L_2(q_1, q_4) = 3 < L_2(q_1, q_7) = 6$$

Is "I" closer to James than to McGuire?

# Memory

o A representation of the past

o Project information at timestep $t$ on a latent space $c_t$ using parameters $\theta$

o Re-using the projected information from $t$ at $t+1$

$$c_{t+1} = h(x_{t+1}, c_t; \theta)$$

o Recurrent parameters $\theta$ are the shared for all timesteps $t = 0, \dots$

$$c_{t+1} = h(x_{t+1}, h(x_t, h(x_{t-1}, \dots h(x_1, c_0; \theta); \theta); \theta); \theta)$$

# Memory as a Graph

o Simplest model
  ◦ Input with parameters $U$
  ◦ Memory embedding with parameters $W$
  ◦ Output with parameters $V$

Output $y_t$

Output parameters $V$

Memory parameters $W$

$c_t$

Memory

Input parameters $U$

Input $x_t$

# Memory as a Graph

o Simplest model
  ◦ Input with parameters $U$
  ◦ Memory embedding with parameters $W$
  ◦ Output with parameters $V$

# Memory as a Graph

○ Simplest model

   ◦ Input with parameters $U$

   ◦ Memory embedding with parameters $W$

   ◦ Output with parameters $V$

# Folding the memory

# Recurrent Neural Network (RNN)

○ Only **two** equations

$$c_t = \tanh(U\ x_t + W\ c_{t-1})$$
$$y_t = \text{softmax}(V\ c_t)$$

# RNN Example

o Vocabulary of 5 words

o A memory of 3 units [Hyperparameter that we choose like layer size]
  ◦ $c_t$: $[3 \times 1]$, W: $[3 \times 3]$

o An input projection of 3 dimensions
  ◦ U: $[3 \times 5]$

o An output projections of 10 dimensions
  ◦ V: $[10 \times 3]$

$$c_t = \tanh(\textcolor{red}{U}\, x_t + \textcolor{cyan}{W}\, c_{t-1})$$

$$y_t = \mathrm{softmax}(\textcolor{green}{V}\, c_t)$$

$$\textcolor{red}{U} \cdot x_{t=4} = \begin{bmatrix} 0.1 & -0.3 & 1.2 & 0.6 & -0.8 \\ -0.2 & 0.4 & 0.5 & 0.9 & -0.1 \\ -0.1 & 0.2 & -0.7 & -0.8 & 0.3 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.6 \\ 0.9 \\ -0.8 \end{bmatrix} = U^{(4)}$$

# Rolled Network vs. Multi-layer Network?

o What is really different?

◦ Steps instead of layers

◦ Step parameters shared whereas in a Multi-Layer Network they are different



3-gram Unrolled Recurrent Network

3-layer Neural Network

# Rolled – Unrolled networks

○ What is really different?
  ◦ Steps instead of layers
  ◦ Step parameters shared whereas in a Multi-Layer Network they are different

- Sometimes intermediate outputs are not even needed
- Removing them, we almost end up to a standard Neural Network



3-gram Unrolled Recurrent Network

3-layer Neural Network

# Training Recurrent Networks

o Cross-entropy loss

$$P = \prod_{t,k} y_{tk}^{l_{tk}} \quad \Rightarrow \quad \mathcal{L} = -\log P = \sum_t \mathcal{L}_t = -\frac{1}{T} \sum_t l_t \log y_t$$

o Backpropagation Through Time (BPTT)

◦ Again, chain rule

◦ Only difference: Gradients survive over time steps

# Backpropagation Through Time: An Example

- $\dfrac{\partial \mathcal{L}}{\partial V}, \dfrac{\partial \mathcal{L}}{\partial W}, \dfrac{\partial \mathcal{L}}{\partial U}$

Step by step explanation at:

http://www.wildml.com/2015/10/recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-gradients/

- To make it simpler let's focus on step 3

$$\frac{\partial \mathcal{L}_3}{\partial V}, \frac{\partial \mathcal{L}_3}{\partial W}, \frac{\partial \mathcal{L}_3}{\partial U}$$

$$c_t = \tanh(U\, x_t + W c_{t-1})$$
$$y_t = \mathrm{softmax}(V c_t)$$
$$\mathcal{L} = -\sum_t l_t \log y_t = \sum_t \mathcal{L}_t$$

# Backpropagation Through Time

$$\frac{\partial \mathcal{L}_3}{\partial V} = \frac{\partial \mathcal{L}_3}{\partial y_3} \frac{\partial y_3}{\partial V} = (y_3 - l_3) \cdot c_3$$



$$c_t = \tanh(U\,x_t + Wc_{t-1})$$
$$y_t = \text{softmax}(Vc_t)$$
$$\mathcal{L} = -\sum_t l_t \log y_t = \sum_t \mathcal{L}_t$$

# Backpropagation Through Time

○ $\dfrac{\partial \mathcal{L}_3}{\partial W} = \dfrac{\partial \mathcal{L}_3}{\partial y_3} \dfrac{\partial y_3}{\partial c_3} \dfrac{\partial c_3}{\partial W}$

○ What is the relation between $c_3$ and $W$?
  ◦ Two-fold: $c_t = \tanh(U\,x_t + W\,c_{t-1})$

○ $\dfrac{\partial\, f(\varphi(x), \psi(x))}{\partial x} = \dfrac{\partial f}{\partial \varphi} \dfrac{\partial \varphi}{\partial x} + \dfrac{\partial f}{\partial \psi} \dfrac{\partial \psi}{\partial x}$

○ $\dfrac{\partial c_3}{\partial W} \propto c_2 + \dfrac{\partial c_2}{\partial W} \quad (\dfrac{\partial W}{\partial W} = 1)$

$$c_t = \tanh(U\,x_t + W c_{t-1})$$
$$y_t = \mathrm{softmax}(V c_t)$$
$$\mathcal{L} = -\sum_t l_t \log y_t = \sum_t \mathcal{L}_t$$

# Recursively

$\circ \quad \dfrac{\partial c_3}{\partial W} = c_2 + \dfrac{\partial c_2}{\partial W}$

$\circ \quad \dfrac{\partial c_2}{\partial W} = c_1 + \dfrac{\partial c_1}{\partial W}$

$\circ \quad \dfrac{\partial c_1}{\partial W} = c_0 + \dfrac{\partial c_0}{\partial W}$

$\dfrac{\partial c_3}{\partial W} = \displaystyle\sum_{t=1}^{3} \dfrac{\partial c_3}{\partial c_t} \dfrac{\partial c_t}{\partial W} \Rightarrow \boxed{\dfrac{\partial \mathcal{L}_3}{\partial W} = \displaystyle\sum_{t=1}^{3} \dfrac{\partial \mathcal{L}_3}{\partial y_3} \dfrac{\partial y_3}{\partial c_3} \dfrac{\partial c_3}{\partial c_t} \dfrac{\partial c_t}{\partial W}}$

$c_t = \tanh(U\,x_t + W c_{t-1})$

$y_t = \text{softmax}(V c_t)$

$\mathcal{L} = -\displaystyle\sum_t l_t \log y_t = \sum_t \mathcal{L}_t$

$y_1, \mathcal{L}_1 \qquad y_2, \mathcal{L}_2 \qquad y_3, \mathcal{L}_3$

$V \qquad V \qquad V$

$W \qquad W \qquad W \qquad W$

$U \qquad U \qquad U$

$x_1 \qquad x_2 \qquad x_3$

# What makes RNNs tick?

o The latent memory space is composed of multiple dimensions

o A subspace of the memory state space can store information if multiple basins ◆ of attraction in some dimensions exist

o Gradients must be strong near the basin boundaries

# Training RNNs is hard

o Vanishing gradients
  ◦ After a few time steps the gradients become almost 0

o Exploding gradients
  ◦ After a few time steps the gradients become huge

o Can't capture long-term dependencies

# Alternative formulation for RNNs

o An alternative formulation to derive conclusions and intuitions

$$c_t = W \cdot \tanh(c_{t-1}) + U \cdot x_t + b$$

$$\mathcal{L} = \sum_t \mathcal{L}_t(c_t)$$

# Another look at the gradients

- $\mathcal{L} = L(c_T(c_{T-1}(\ldots(c_1(x_1, c_0; W); W); W); W)$

- $\dfrac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \dfrac{\partial \mathcal{L}_t}{\partial c_t} \dfrac{\partial c_t}{\partial c_\tau} \dfrac{\partial c_\tau}{\partial W}$

- $\dfrac{\partial \mathcal{L}}{\partial c_t} \dfrac{\partial c_t}{\partial c_\tau} = \dfrac{\partial \mathcal{L}}{\partial c_t} \cdot \dfrac{\partial c_t}{\partial c_{t-1}} \cdot \dfrac{\partial c_{t-1}}{\partial c_{t-2}} \cdot \ldots \cdot \dfrac{\partial c_{\tau+1}}{\partial c_\tau} \leq \eta^{t-\tau} \dfrac{\partial \mathcal{L}_t}{\partial c_t}$

  $\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$
  $Rest \rightarrow short\text{-}term\ factors$

  $\underbrace{\phantom{xxxxxxxxx}}$
  $t \gg \tau \rightarrow long\text{-}term\ factors$

  $\eta$ determines the norm of the gradients

- The RNN gradient is a recursive product of $\dfrac{\partial c_t}{\partial c_{t-1}}$

# RNN gradients in 1D

○ $\dfrac{\partial \mathcal{L}}{\partial c_t} = \dfrac{\partial \mathcal{L}}{\partial c_T} \cdot \underbrace{\dfrac{\partial c_T}{\partial c_{T-1}}}_{<\,1} \cdot \underbrace{\dfrac{\partial c_{T-1}}{\partial c_{T-2}}}_{<\,1} \cdot \; \dots \; \cdot \underbrace{\dfrac{\partial c_{t+1}}{\partial c_{c_t}}}_{<\,1} \Bigg\}$ $\quad \dfrac{\partial \mathcal{L}}{\partial W} \ll 1 \Longrightarrow$ Vanishing gradient

○ $\dfrac{\partial \mathcal{L}}{\partial c_t} = \dfrac{\partial \mathcal{L}}{\partial c_T} \cdot \underbrace{\dfrac{\partial c_T}{\partial c_{T-1}}}_{>\,1} \cdot \underbrace{\dfrac{\partial c_{T-1}}{\partial c_{T-2}}}_{>\,1} \cdot \; \dots \; \cdot \underbrace{\dfrac{\partial c_1}{\partial c_{c_t}}}_{>\,1} \Bigg\}$ $\quad \dfrac{\partial \mathcal{L}}{\partial W} \gg 1 \Longrightarrow$ Exploding gradient

# RNN gradients in N-D

o When $c_T \in \mathbb{R}^N$ then $\frac{\partial c_t}{\partial c_{t-1}}$ is a Jacobian

o $\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \underbrace{\frac{\partial c_T}{\partial c_{T-1}}}_{< 1} \cdot \underbrace{\frac{\partial c_{T-1}}{\partial c_{T-2}}}_{< 1} \cdot \ldots \cdot \underbrace{\frac{\partial c_{t+1}}{\partial c_t}}_{< 1} \Bigg\}$     $\frac{\partial \mathcal{L}}{\partial \theta} \ll 1 \Longrightarrow$ Vanishing gradient

o $\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \underbrace{\frac{\partial c_T}{\partial c_{T-1}}}_{> 1} \cdot \underbrace{\frac{\partial c_{T-1}}{\partial c_{T-2}}}_{> 1} \cdot \ldots \cdot \underbrace{\frac{\partial c_{t+1}}{\partial c_t}}_{> 1} \Bigg\}$     $\frac{\partial \mathcal{L}}{\partial \theta} \gg 1 \Longrightarrow$ Exploding gradient

# The Jacobian

$$y \in \mathbb{R}^2, x \in \mathbb{R}^3 : \quad \frac{d\boldsymbol{y}}{d\boldsymbol{x}} = \begin{bmatrix} \dfrac{\partial y^{(1)}}{\partial x^{(1)}} & \dfrac{\partial y^{(1)}}{\partial x^{(2)}} & \dfrac{\partial y^{(1)}}{\partial x^{(3)}} \\ \dfrac{\partial y^{(2)}}{\partial x^{(1)}} & \dfrac{\partial y^{(2)}}{\partial x^{(2)}} & \dfrac{\partial y^{(2)}}{\partial x^{(3)}} \end{bmatrix}$$

# RNN gradients in N-D

- When $c_T \in \mathbb{R}^N$ then $\frac{\partial c_t}{\partial c_{t-1}}$ is a Jacobian

- Spectral radius ρ (~largest eigenvalue) of Jacobian is important

- $\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \ldots \cdot \frac{\partial c_{t+1}}{\partial c_{c_t}}$
  $\rho < 1 \quad\quad \rho < 1 \quad\quad\quad \rho < 1$
  $\Bigg\}$ $\quad \frac{\partial \mathcal{L}}{\partial c_t} \ll 1 \Longrightarrow$ **Vanishing gradient**

- $\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \ldots \cdot \frac{\partial c_{t+1}}{\partial c_t}$
  $\rho > 1 \quad\quad \rho > 1 \quad\quad\quad \rho > 1$
  $\Bigg\}$ $\quad \frac{\partial \mathcal{L}}{\partial c_t} \gg 1 \Longrightarrow$ **Exploding gradient**

# Gradient clipping for exploding gradients

○ Scale the gradients to a threshold



**Pseudocode**

1. $g \leftarrow \frac{\partial \mathcal{L}}{\partial W}$

2. $\text{if } \|g\| > \theta_0:$

$$g \leftarrow \frac{\theta_0}{\|g\|} g$$

$\text{else:}$

$\text{print('Do nothing')}$

○ Simple, but works!

# Vanishing gradients

o The gradient of the error w.r.t. to intermediate cell

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial c_t}{\partial c_\tau} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

# Vanishing gradients

- For $t = 1, r = 2 \implies \dfrac{\partial \mathcal{L}_2}{\partial W} \propto \dfrac{\partial c_2}{\partial c_1}$

- For $t = 1, r = 3 \implies \dfrac{\partial \mathcal{L}_3}{\partial W} \propto \dfrac{\partial c_3}{\partial c_1} = \dfrac{\partial c_3}{\partial c_2} \cdot \dfrac{\partial c_2}{\partial c_1}$

- For $t = 1, r = 4 \implies \dfrac{\partial \mathcal{L}_4}{\partial W} \propto \dfrac{\partial c_4}{\partial c_1} = \dfrac{\partial c_4}{\partial c_3} \cdot \dfrac{\partial c_3}{\partial c_2} \cdot \dfrac{\partial c_2}{\partial c_1}$

# Vanishing gradients

o The gradient of the error w.r.t. to intermediate cell

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial c_t}{\partial c_\tau} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

o Long-term dependencies get exponentially smaller weights

# Rescaling vanishing gradients?

o Not good solution

o Weights are shared between timesteps → Loss summed over timesteps

$$\mathcal{L} = \sum_t \mathcal{L}_t \quad \Rightarrow \quad \frac{\partial \mathcal{L}}{\partial W} = \sum_t \frac{\partial \mathcal{L}_t}{\partial W}$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W} = \sum_{\tau=1}^t \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \textcolor{red}{\frac{\partial c_\tau}{\partial W}}$$

o Rescaling for one timestep ($\frac{\partial \mathcal{L}_t}{\partial W}$) affects all timesteps

◦ The rescaling factor for one timestep does not work for another

# More intuitively



$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}_1}{\partial W} + \frac{\partial \mathcal{L}_2}{\partial W} + \frac{\partial \mathcal{L}_3}{\partial W} + \frac{\partial \mathcal{L}_4}{\partial W} + \frac{\partial \mathcal{L}_5}{\partial W}$$

$\frac{\partial \mathcal{L}_1}{\partial W}$

$\frac{\partial \mathcal{L}_2}{\partial W}$

$\frac{\partial \mathcal{L}_3}{\partial W}$

$\frac{\partial \mathcal{L}_4}{\partial W}$

$\frac{\partial \mathcal{L}_5}{\partial W}$

# More intuitively

- Let's say $\frac{\partial \mathcal{L}_1}{\partial w} \propto 1, \frac{\partial \mathcal{L}_2}{\partial w} \propto 1/10, \frac{\partial \mathcal{L}_3}{\partial w} \propto 1/100, \frac{\partial \mathcal{L}_4}{\partial w} \propto 1/1000, \frac{\partial \mathcal{L}_5}{\partial w} \propto 1/10000$

- $\frac{\partial \mathcal{L}}{\partial w} = \sum_r \frac{\partial \mathcal{L}_r}{\partial w} = 1.1111$

- If $\frac{\partial \mathcal{L}}{\partial w}$ rescaled to 1 $\rightarrow \frac{\partial \mathcal{L}_5}{\partial w} \propto 10^{-5}$

- Longer-term dependencies negligible
  - Weak recurrent modelling
  - Learning focuses on the short-term only

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}_1}{\partial w} + \frac{\partial \mathcal{L}_2}{\partial w} + \frac{\partial \mathcal{L}_3}{\partial w} + \frac{\partial \mathcal{L}_4}{\partial w} + \frac{\partial \mathcal{L}_5}{\partial w}$$

# Recurrent networks ∝ Dynamical systems

○ In the figures $x_t \propto c_t$ and $x_t \propto F(Wx_{t-1} + Uu_t + b)$



*Figure 4.* This diagram illustrates how the change in $\mathbf{x}_t$, $\Delta\mathbf{x}_t$, can be large for a small $\Delta\mathbf{x}_0$. The blue vs red (left vs right) trajectories are generated by the same maps $F_1, F_2, \ldots$ for two different initial states.



*Figure 5.* Illustrates how one can break apart the maps $F_1, ..F_t$ into a constant map $\tilde{F}$ and the maps $U_1, .., U_t$. The dotted vertical line represents the boundary between basins of attraction, and the straight dashed arrow the direction of the map $\tilde{F}$ on each side of the boundary. This diagram is an extension of Fig. 4.

# Fixing vanishing gradients

o Regularization on the recurrent weights
  ◦ Force the error signal not to vanish

$$\Omega = \sum_t \Omega_t = \sum_t \left( \frac{\left| \frac{\partial \mathcal{L}}{\partial c_{t+1}} \frac{\partial c_{t+1}}{\partial c_t} \right|}{\left| \frac{\partial \mathcal{L}}{\partial c_{t+1}} \right|} - 1 \right)^2$$

o Advanced recurrent modules

o Long-Short Term Memory module

o Gated Recurrent Unit module

Pascanu et al., On the diculty of training Recurrent Neural Networks, 2013

# Advanced RNNs

# How to fix the vanishing gradients?

- Error signal over time must have not too large, not too small norm

- Solution: have an activation function with derivative equal to 1
  - Identify function

- By doing so, gradients do not become too small not too large

# Long Short-Term Memory (LSTM: Beefed up RNN)

$$i = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$

$$f = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$

$$o = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$

$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$

$$m_t = \tanh(c_t) \odot o$$

More info at:

# Cell state

o The cell state carries the essential information over time

Cell state line

$$i = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

$$f = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

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

$$m_t = \tanh(c_t) \odot o$$

# LSTM nonlinearities

○ $\sigma \in (0, 1)$: control gate – something like a switch

○ tanh $\in (-1, 1)$: recurrent nonlinearity

$$i = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$

$$f = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$

$$o = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$

$$\tilde{c}_t = \tanh\big(x_t U^{(g)} + m_{t-1} W^{(g)}\big)$$

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

$$m_t = \tanh(c_t) \odot o$$

# LSTM Step-by-Step: Step (1)

○ E.g. Model the sentence "Yesterday she slapped me. Today she loves me."

○ Decide what to forget and what to remember for the new memory

◦ Sigmoid 1 → Remember everything

◦ Sigmoid 0 → Forget everything

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\tilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

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

$$m_t = \tanh(c_t) \odot o$$

# LSTM Step-by-Step: Step (2)

○ Decide what new information should you add to the new memory

  ◦ Modulate the input $i_t$

  ◦ Generate candidate memories $\widetilde{c}_t$

$$i_t = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$

$$f_t = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$

$$o_t = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$

$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$

$$m_t = \tanh(c_t) \odot o$$

# LSTM Step-by-Step: Step (3)

○ Compute and update the current cell state $c_t$

  ◦ Depends on the previous cell state

  ◦ What we decide to forget

  ◦ What inputs we allow

  ◦ The candidate memories

$$i_t = \sigma\left(x_t U^{(i)} + m_{t-1} W^{(i)}\right)$$

$$f_t = \sigma\left(x_t U^{(f)} + m_{t-1} W^{(f)}\right)$$

$$o_t = \sigma\left(x_t U^{(o)} + m_{t-1} W^{(o)}\right)$$

$$\widetilde{c}_t = \tanh\left(x_t U^{(g)} + m_{t-1} W^{(g)}\right)$$

$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$

$$m_t = \tanh(c_t) \odot o$$

# LSTM Step-by-Step: Step (4)

o Modulate the output
  ◦ Does the cell state contain something relevant? → Sigmoid 1

o Generate the new memory

$$i_t = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$$

$$f_t = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$$

$$o_t = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$$

$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$

$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$

$$m_t = \tanh(c_t) \odot o$$

# LSTM Unrolled Network

○ Macroscopically very similar to standard RNNs

○ The engine is a bit different (more complicated)
  ◦ Because of their gates LSTMs capture long and short term dependencies

# Beyond RNN & LSTM

o LSTM with peephole connections
   ◦ Gates have access also to the previous cell states $c_{t-1}$ (not only memories)
   ◦ Coupled forget and input gates, $c_t = f_t \odot c_{t-1} + (1 - f_t) \odot \widetilde{c}_t$
   ◦ Bi-directional recurrent networks

o Gated Recurrent Units (GRU)

o Deep recurrent architectures

o Recursive neural networks
   ◦ Tree structured

o Multiplicative interactions

o Generative recurrent architectures

| LSTM (2) | → | LSTM (2) | → | LSTM (2) | → |
| LSTM (1) | → | LSTM (1) | → | LSTM (1) | → |

# Applications of Recurrent Networks

a man in a suit and tie standing in front of a building

NeuralTalk and Walk, recognition, text description of the image while walking

CloudCV: Visual Question Answering (VQA)

More details about the VQA dataset can be found here.
State-of-the-art VQA model and code available here

CloudCV can answer questions you ask about an image

Browsers currently supported: Google Chrome, Mozilla Firefox

Try CloudCV VQA: Sample Images

Click on one of these images to send it to our servers (Or upload your own images below)

55 minutes to work
Light traffic on 101

Embarcadero
Train station

Hi Motherboard readers!

This entire post was hand written by a neural network.

(It probably writes better than you.)

Of course, a neural network doesn't actually have hands

And the original text was typed by me, a human.

So what's going on here?

A neural network is a program that can learn to follow a set of rules

But it can't do it alone. It needs to be trained.

This neural network was trained on a corpus of writing samples.

...but of the locations of a pen-tip as people write.

This is how the network learns and creates different styles from prior examples.

And it can use this knowledge to generate handwritten notes from inputted text.

It can create its own style, or mimic another's.

No two notes are the same.

It's the work of Alex Graves at the University of Toronto

And you can try it too!

# Machine Translation

- The phrase in the source language is one sequence
  - "Today the weather is good"

- The phrase in the target language is also a sequence
  - "Vandaag is het weer goed"

- Problems
  - no perfect word alignment, sentence length might differ

- Solution
  - Encoder-decoder scheme

# Better Machine Translation

- It might even pay off reversing the source sentence
  - The first target words will be closer to their respective source words

- The encoder and decoder parts can be modelled with different LSTMs

- Deep LSTM

# Image captioning

o An image is a thousand words, literally!

o Pretty much the same as machine transation

o Replace the encoder part with the output of a Convnet
  ◦ E.g. use Alexnet or a VGG16 network

o Keep the decoder part to operate as a translator

# Question answering

- Bleeding-edge research, no real consensus
  - Very interesting open, research problems

- Again, pretty much like machine translation

- Again, Encoder-Decoder paradigm
  - Insert the question to the encoder part
  - Model the answer at the decoder part

- Question answering with images also
  - Again, bleeding-edge research
  - How/where to add the image?
  - What has been working so far is to add the image only in the beginning

Q: John entered the living room, where he met Mary. She was drinking some wine and watching a movie. What room did John enter?

A: John entered the living room.



Q: what are the people playing?

A: They play beach football

# Summary

o Recurrent Neural Networks (RNN) for sequences

o Backpropagation Through Time

o Vanishing and Exploding Gradients and Remedies

o RNNs using Long Short-Term Memory (LSTM)

o Applications of Recurrent Neural Networks

# Reading material & references

- [http://www.deeplearningbook.org/](http://www.deeplearningbook.org/)
  - Part II: Chapter 10

- Excellent blog post on Backpropagation Through Time
  - [http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/](http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/)
  - [http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-2-implementing-a-language-model-rnn-with-python-numpy-and-theano/](http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-2-implementing-a-language-model-rnn-with-python-numpy-and-theano/)
  - [http://www.wildml.com/2015/10/recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-gradients/](http://www.wildml.com/2015/10/recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-gradients/)

- Excellent blog post explaining LSTMs
  - [http://colah.github.io/posts/2015-08-Understanding-LSTMs/](http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

[Pascanu2013] Pascanu, Mikolov, Bengio. On the difficulty of training Recurrent Neural Networks, JMLR, 2013

# Next lecture

- Memory networks
- Recursive networks