Recurrent Neural Networks (part 2 -- technical things)

Yoav Goldberg

Concrete RNN Functions







enc(what is your name)

- Very strong models of sequential data.
- **Trainable** function from *n* vectors to a single vector.



- There are different variants (implementations).
- So far, we focused on the interface level.



- There's a vector \mathbf{y}_i for every prefix $\mathbf{x}_{1:i}$

 $\mathbf{X_{i}}$

• What do the RNN function look like?

 θ



But we can train them.
 define function form define loss

$$R_{CBOW}(\mathbf{s_{i-1}}, \mathbf{x_i}) = \mathbf{s_{i-1}} + \mathbf{x_i}$$

(what are the parameters?)

$$R_{CBOW}(\mathbf{s_{i-1}}, \mathbf{x_i}) = \mathbf{s_{i-1}} + \mathbf{x_i}$$

(what are the parameters?)

$$R_{CBOW}(\mathbf{s_{i-1}}, x_i) = \mathbf{s_{i-1}} + \mathbf{E}_{[x_i]}$$

Is this a good parameterization?

$$R_{CBOW}(\mathbf{s_{i-1}}, x_i) = \mathbf{s_{i-1}} + \mathbf{E}_{[x_i]}$$

how about this modification?

$R_{CBOW}(\mathbf{s_{i-1}}, x_i) = \underline{\tanh}(\mathbf{s_{i-1}} + \mathbf{E}_{[x_i]})$

 $R_{SRNN}(\mathbf{s_{i-1}}, \mathbf{x_i}) = tanh(\mathbf{W^s} \cdot \mathbf{s_{i-1}} + \mathbf{W^x} \cdot \mathbf{x_i})$

$$R_{SRNN}(\mathbf{s_{i-1}}, \mathbf{x_i}) = tanh(\mathbf{W^s} \cdot \mathbf{s_{i-1}} + \mathbf{W^x} \cdot \mathbf{x_i})$$

- Looks very simple.
- Theoretically very powerful.
- In practice not so much (hard to train).
- Why? Vanishing gradients.

$$R_{SRNN}(\mathbf{s_{i-1}}, \mathbf{x_i}) = tanh(\mathbf{W^s} \cdot \mathbf{s_{i-1}} + \mathbf{W^x} \cdot \mathbf{x_i})$$

Another view on behavior:

- RNN as a "computer": input **xi** arrives, memory **s** is updated.
- In the Elman RNN, **entire memory is written** at each time-step.



Another view on behavior:

- RNN as a "computer": input **xi** arrives, memory **s** is updated.
- In the Elman RNN, entire memory is written at each time-step.

LSTM RNN

$$\begin{aligned} R_{LSTM}(\mathbf{s_{j-1}}, \mathbf{x_j}) = [\mathbf{c_j}; \mathbf{h_j}] \\ \mathbf{c_j} = \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i} \\ \mathbf{h_j} = \tanh(\mathbf{c_j}) \odot \mathbf{o} \\ \mathbf{i} = \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}}) \\ \mathbf{f} = \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}}) \\ \mathbf{o} = \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}}) \\ \mathbf{g} = \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}}) \\ O_{LSTM}(\mathbf{s_j}) = O_{LSTM}([\mathbf{c_j}; \mathbf{h_j}]) = \mathbf{h_j} \end{aligned}$$

LSTM RNN

better controlled memory access

continuous gates

Differentiable "Gates"

- The main idea behind the LSTM is that you want to somehow control the "memory access".
- In a SimpleRNN:

$$R_{SRNN}(\mathbf{s_{i-1}}, \mathbf{x_i}) = tanh(\mathbf{W^s} \cdot \mathbf{s_{i-1}} + \mathbf{W^x} \cdot \mathbf{x_i})$$

read previous state memory write new input

• All the memory gets overwritten

- We'd like to:
 - * Selectively read from some memory "cells".
 - * Selectively write to some memory "cells".

- We'd like to:
 - * Selectively read from some memory "cells".
 - * Selectively write to some memory "cells".



- We'd like to:
 - * Selectively read from some memory "cells".
 - * Selectively write to some memory "cells".
- A gate function:



• Using the gate function to control access:



• Using the gate function to control access:



• (can also tie them: $\mathbf{g}^{\mathbf{r}} = 1 - \mathbf{g}^{\mathbf{w}}$)



Differentiable "Gates"

Problem with the gates:

- * they are fixed.
- * they don't depend on the input or the output.

Differentiable "Gates"

Problem with the gates:

- * they are fixed.
- * they don't depend on the input or the output.
- Solution: make them smooth, input dependent, and trainable.



LSTM (Long short-term Memory)

• The LSTM is a specific combination of gates.

$$R_{LSTM}(\mathbf{s_{j-1}}, \mathbf{x_j}) = [\mathbf{c_j}; \mathbf{h_j}]$$

$$\mathbf{c_j} = \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h_j} = \tanh(\mathbf{c_j}) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{f} = \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{o} = \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{g} = \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}})$$

$$O_{LSTM}(\mathbf{s_j}) = O_{LSTM}([\mathbf{c_j}; \mathbf{h_j}]) = \mathbf{h_j}$$

• The GRU is a different combination of gates.

$$\mathbf{s_j} = R_{\text{GRU}}(\mathbf{s_{j-1}}, \mathbf{x_j}) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s_{j-1}} + \mathbf{z} \odot \tilde{\mathbf{s_j}}$$
$$\mathbf{z} = \sigma(\mathbf{x_j} \mathbf{W^{xz}} + \mathbf{s_{j-1}} \mathbf{W^{sz}})$$
$$\mathbf{r} = \sigma(\mathbf{x_j} \mathbf{W^{xr}} + \mathbf{s_{j-1}} \mathbf{W^{sr}})$$
$$\tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j} \mathbf{W^{xs}} + (\mathbf{r} \odot \mathbf{s_{j-1}}) \mathbf{W^{sg}})$$

GRU vs LSTM

- The GRU and the LSTM are very similar ideas.
- Invented independently of the LSTM, almost two decades later.

• The GRU formulation:

$$\mathbf{s_j} = R_{\mathrm{GRU}}(\mathbf{s_{j-1}}, \mathbf{x_j}) =$$

 $\label{eq:proposal} \mbox{Proposal state:} \qquad \ \ \tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j}\mathbf{W^{xs}} + (\mathbf{r}\odot\mathbf{s_{j-1}})\mathbf{W^{sg}})$

• The GRU formulation:

$$\mathbf{s_j} = R_{\mathrm{GRU}}(\mathbf{s_{j-1}}, \mathbf{x_j}) =$$

gate controlling effect of prev on proposal:

$$\mathbf{r} = \sigma(\mathbf{x_j}\mathbf{W^{xr}} + \mathbf{s_{j-1}}\mathbf{W^{sr}})$$
$$\tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j}\mathbf{W^{xs}} + \mathbf{(r)} \cdot \mathbf{s_{j-1}})\mathbf{W^{sg}})$$

$\begin{array}{l} \text{blend of old state and} \\ \textbf{proposal state} \\ \textbf{s_j} = R_{\text{GRU}}(\textbf{s_{j-1}}, \textbf{x_j}) = (\textbf{1} - \textbf{z}) \odot \textbf{s_{j-1}} + \textbf{z} \odot \tilde{\textbf{s_j}} \end{array}$

$$\mathbf{r} = \sigma(\mathbf{x_j}\mathbf{W^{xr}} + \mathbf{s_{j-1}}\mathbf{W^{sr}})$$
$$\tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j}\mathbf{W^{xs}} + (\mathbf{r} \odot \mathbf{s_{j-1}})\mathbf{W^{sg}})$$

$$\begin{split} \mathbf{s_j} &= R_{\mathrm{GRU}}(\mathbf{s_{j-1}}, \mathbf{x_j}) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s_{j-1}} + \mathbf{z} \odot \tilde{\mathbf{s_j}} \\ \\ \textbf{gate for controlling} & \mathbf{z} = \sigma(\mathbf{x_j} \mathbf{W^{xz}} + \mathbf{s_{j-1}} \mathbf{W^{sz}}) \\ \textbf{the blend} & \mathbf{r} = \sigma(\mathbf{x_j} \mathbf{W^{xr}} + \mathbf{s_{j-1}} \mathbf{W^{sr}}) \\ & \tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j} \mathbf{W^{xs}} + (\mathbf{r} \odot \mathbf{s_{j-1}}) \mathbf{W^{sg}}) \end{split}$$

$$\begin{aligned} \mathbf{s_j} &= R_{\mathrm{GRU}}(\mathbf{s_{j-1}}, \mathbf{x_j}) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s_{j-1}} + \mathbf{z} \odot \tilde{\mathbf{s_j}} \\ & \mathbf{z} = \sigma(\mathbf{x_j} \mathbf{W^{xz}} + \mathbf{s_{j-1}} \mathbf{W^{sz}}) \\ & \mathbf{r} = \sigma(\mathbf{x_j} \mathbf{W^{xr}} + \mathbf{s_{j-1}} \mathbf{W^{sr}}) \\ & \tilde{\mathbf{s_j}} = \tanh(\mathbf{x_j} \mathbf{W^{xs}} + (\mathbf{r} \odot \mathbf{s_{j-1}}) \mathbf{W^{sg}}) \end{aligned}$$

LSTM (Long short-term Memory)

• The LSTM is formulation:

 $\begin{aligned} R_{LSTM}(\mathbf{s_{j-1}},\mathbf{x_j}) = & [\mathbf{c_j};\mathbf{h_j}] \\ \mathbf{c_j} = & \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i} \\ \mathbf{h_j} = & \tanh(\mathbf{c_j}) \odot \mathbf{o} \\ & \mathbf{i} = & \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}}) \\ & \mathbf{f} = & \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}}) \\ & \mathbf{o} = & \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}}) \\ & \mathbf{g} = & \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}}) \end{aligned}$

$$O_{LSTM}(\mathbf{s_j}) = O_{LSTM}([\mathbf{c_j}; \mathbf{h_j}]) = \mathbf{h_j}$$

LSTM (Long short-term Memory)

• The LSTM is formulation:

$$\begin{split} R_{LSTM}(\mathbf{s_{j-1}},\mathbf{x_j}) =& [\mathbf{c_j};\mathbf{h_j}] \\ & \mathbf{c_j} =& \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i} \\ & \mathbf{h_j} =& \tanh(\mathbf{c_j}) \odot \mathbf{o} \\ \text{"input"} & \mathbf{i} =& \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}}) \\ \text{"forget"} & \mathbf{f} =& \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}}) \\ \text{"output"} & \mathbf{o} =& \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}}) \\ & \mathbf{g} =& \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}}) \\ \end{split}$$

 $O_{LSTM}(\mathbf{s_j}) = O_{LSTM}([\mathbf{c_j}; \mathbf{h_j}]) = \mathbf{h_j}$
Read More

- The gated architecture also helps the vanishing gradients problems.
- For a good explanation, see Kyunghyun Cho's notes: <u>http://arxiv.org/abs/1511.07916</u> sections 4.2, 4.3
- Chris Olah's blog post (link in class website)

Other Variants

- Many other variants exist.
- Mostly perform similarly to each other.
 - Different tasks may work better with different variants.
- The important idea is the differentiable gates.

Differences?

- There ARE formal difference in power between the GRU and the LSTM.
 - The LSTM can count, the GRU cannot. See paper by Weiss, Goldberg and Yahav in reading material. <u>https://arxiv.org/abs/1805.04908</u>

- Still an open question how to perform well.
- One suggestion:



each color is a different random dropout mask

- Still an open question how to perform well.
- Yarin Gal's Dropout:



each color is a different random dropout mask

- Still an open question how to perform well.
- Yarin Gal's Dropout:



each color is a different random dropout mask

Practicalities

- Most toolkits require a **fixed** computation graph for all examples.
- But RNNs have different input lengths. What do we do?
 - Option 1: Use a tool that does not pose this limitation.
 - Option 2: Enforce max length + 0 padding for shorter sequences.







Transducers for POS tagging?

Can predict the tag of word i based on words 1,...,i-1, but...





Each state encodes the entire history up to that state. This is not bad. But what about the future?

















X_{fox}

Xjumped

 \mathbf{X}_{*}

One RNN runs left to right. Another runs right to left. Encode **both future and history** of a word.

Xbrown

 $\mathbf{x_{the}}$



 $BiRNN(\mathbf{x_{1:7}}, 4) = [\mathbf{y_4^F}; \mathbf{y_4^R}]$ $\mathbf{y_4^F} = RNN_F(\mathbf{x_{1:4}})$ $\mathbf{y_4^R} = RNN_R(\mathbf{x_{7:4}})$

Encode **both future and history** of a word.

One RNN runs left to right.

Another runs right to left.



"deep RNNs"



RNN can be stacked deeper is better! (better how?)



Deep Bi-RNNs



BI-RNN can also be stacked



(Deep) BI-RNNs

- provide an "infinite" window around a focus word.
- learn to extract what's important.
- easy to train!
- very effective for sequence tagging.

Sequence tagger (here, POS)



Sequence tagger (here, POS)



RNNs can be easily "nested". Output of one RNN feeds into another RNN.

Sequence tagger (here, POS)



Back-off to char-level RNNs for Unknown Words:

(inspired by Ling et al, Ballesteros, Dyer and Smith)





Back-off to char-level RNNs for Unknown Words:

(inspired by Ling et al, Ballesteros, Dyer and Smith)



But RNNs need plenty of data to train on!

Do they, really?

— bi-LSTM TnT CRF



Multilingual Part-of-Speech Tagging with Bidirectional Long Short-Term Memory Models and Auxiliary Loss

Barbara Plank University of Groningen The Netherlands b.plank@rug.nl

Anders Søgaard University of Copenhagen Denmark soegaard@hum.ku.dk yoav.goldberg@gmail.com

Yoav Goldberg Bar-Ilan University Israel

But RNNs need plenty of data to train on!

Do they, really?









Germanic



Romance



Semitic



But RNNs need plenty of data to train on! Do they, really? bi-LSTM TnT CRF Indoeuropean non-Indoeuropean Slavic accuracy For POS-tagging, needs less data than CRF, breaks even with HMM at ~500 sentences accuracy sentences sentences sentences



RNN Recap

- Trainable encoders of sequential data.
- Allow capturing non-markovian (Infinite) history.
- Used as "lego bricks", feed into other components.
- Very effective for language modeling.
- Very effective for feature extraction.



BI-RNN's Recap

- Represent the history up to a point in the sequence, and the future from a point in a sequence.
- Feed into an MLP (or linear classifier) to classify the point based on history and future.
- The network learns what's important in the history and future for the given prediction task.
- "Infinite window"

Batching Reminder









- Sequential in nature, very little parallelism.
- (Compare, e.g., to a Convolutional Network)



in transduction mode, can batch the final output predictions
Batching in RNNs

proposals:

newer architectures, that try to make the sequencedependent parts cheaper.

QUASI-RECURRENT NEURAL NETWORKS

James Bradbury; Stephen Merity; Caiming Xiong & Richard Socher Salesforce Research Palo Alto, California {james.bradbury,smerity,cxiong,rsocher}@salesforce.com

Batching in RNNs

proposals:

newer architectures, that try to make the sequencedependent parts cheaper.

QUASI-RECURRENT NEURAL NETWORKS

James Bradbury; Stephen Merity; Caiming Xiong & Richard Socher Salesforce Research Palo Alto, California {james.bradbury,smerity,cxiong,rsocher}@salesforce.com

Attention Is All You Need

("transformers")

Batching in RNNs

Can batch across sequences.













what if the sequences are different lengths?















this is how its done in TF, PyTorch. supported also in DyNet, but...





really annoying super-confusing with biLSTMs practically impossible with our complex networks

Auto Batching



what if the sequences are different lengths?







then create a separate network for each (easy)









nodes in blue are ready to be executed



nodes in red will be executed using batch operations




































note: batching operations, not inputs.

Recent Efficiency Benchmarks

For some tasks, GPU is much better. For other, CPU still wins. (why, how?)

- What can be batched?
- What can be done in efficient math operations?

$$R_{LSTM}(\mathbf{s_{j-1}}, \mathbf{x_j}) = [\mathbf{c_j}; \mathbf{h_j}]$$

$$\mathbf{c_j} = \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h_j} = \tanh(\mathbf{c_j}) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{f} = \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{o} = \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{g} = \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}})$$

$$R_{LSTM}(\mathbf{s_{j-1}}, \mathbf{x_j}) = [\mathbf{c_j}; \mathbf{h_j}]$$

$$\mathbf{c_j} = \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h_j} = \tanh(\mathbf{c_j}) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{f} = \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{o} = \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{g} = \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}})$$

in transduction mode, can batch all inputs (not states)

$$R_{LSTM}(\mathbf{s_{j-1}}, \mathbf{x_j}) = [\mathbf{c_j}; \mathbf{h_j}]$$

$$\mathbf{c_j} = \mathbf{c_{j-1}} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h_j} = \tanh(\mathbf{c_j}) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W^{xi}} \cdot \mathbf{x_j} + \mathbf{W^{hi}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{f} = \sigma(\mathbf{W^{xf}} \cdot \mathbf{x_j} + \mathbf{W^{hf}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{o} = \sigma(\mathbf{W^{xo}} \cdot \mathbf{x_j} + \mathbf{W^{ho}} \cdot \mathbf{h_{j-1}})$$

$$\mathbf{g} = \tanh(\mathbf{W^{xg}} \cdot \mathbf{x_j} + \mathbf{W^{hg}} \cdot \mathbf{h_{j-1}})$$

all gates computations can be done in single mat-mat op.