Advanced Recurrent Architectures

Davide Bacciu

Dipartimento di Informatica Università di Pisa

Intelligent Systems for Pattern Recognition (ISPR)



Outline Motivations

Lecture Outline

- Dealing with structured/compound data
 - Sequence-to-sequence
 - Attention models
- Dealing with very long-term dependencies
 - Multiscale networks
 - Adding memory components
- Neural reasoning
 - Neural Turing machines

Outline Motivations

Gated RNN Refresher





LSTM Cell

GRU Cell

Outline Motivations

Graphical Notation for Compositionality



Outline Motivations

Basic Gated RNN (GRNN) Limitations

- GRNN are excellent to handle size/topology varying data in input
 - How can we handle size/topology varying outputs?
 - Sequence-to-sequence
- Structured data is compound information
 - Efficient processing needs the ability to focus on certain parts of such information
 - Attention mechanism
- GRNN have troubles dealing with very long-range dependencies
 - Introduce multiscale representation explicitly in the architecture
 - Introduce external memory components

Learning to Output Variable Length Sequences

The idea of an unfolded RNN with blank inputs-outputs does not really work well

Output sequence



input sequence

The approach is based on an encoder-decoder scheme



Encoder

Produce a compressed and fixed length representation *c* of all the input sequence



Introduction Structu Advanced RNN Extendi Wrap-up Neural

Structured Output and Attention Extending Memory Neural Reasoning

Decoder



	Structured Output and Attention
Advanced RNN	Extending Memory
	Neural Reasoning

Decoder

С



We risk to lose memory of *c* soon

If we share the parameters between encoder and decoder we can take $s_1 = c$

Or, at least, assume c and s_1 have compatible size

	Structured Output and Attention
Advanced RNN	Extending Memory
	Neural Reasoning

Decoder



c is contextual information kept throughout output generation

	Structured Output and Attention
Advanced RNN	Extending Memory
	Neural Reasoning

Decoder



Structured Output and Attention Extending Memory Neural Reasoning

Sequence-To-Sequence Learning

Encoder-Decoder can share parameters (but it is uncommon)

> Encoder-Decoder can be trained end-to-end or independently





Reversing the input sequence in encoding can increase performance

Structured Output and Attention Extending Memory Neural Reasoning

On the Need of Paying Attention



- Encoder-Decoder scheme assumes the hidden activation of the last input element summarizes sufficient information to generate the output
 - Bias toward most recent past
- Other parts of the input sequence might be very informative for the task
 - Possibly elements appearing very far from sequence end

Structured Output and Attention Extending Memory Neural Reasoning

On the Need of Paying Attention



• Attention mechanism select which part of the sequence to focus on to obtain a good *c*



Structured Output and Attention Extending Memory Neural Reasoning

What's inside of the box?

The Revenge of the Gates!



Structured Output and Attention Extending Memory Neural Reasoning

Opening the Box



Structured Output and Attention Extending Memory Neural Reasoning

Opening the Box – Relevance



Opening the Box – Softmax

Structured Output and Attention Extending Memory Neural Reasoning

Opening the Box – Voting

Structured Output and Attention Extending Memory Neural Reasoning

Attention in Seq2Seq

Structured Output and Attention Extending Memory Neural Reasoning

Learning to Translate with Attention

Bahdanau et al, Show, Neural machine translation by jointly learning to align and translate, ICLR 2015

Structured Output and Attention Extending Memory Neural Reasoning

Seq-to-Seq on Steroids

Yonghui Wu et al. "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation",

Advanced Attention – Generalize Relevance

This component determines how much each h is correlated/associated with current context s

Advanced Attention – Hard Attention

Sample a single encoding using probability α_i

Structured Output and Attention Extending Memory Neural Reasoning

Advanced Attention – Self Attention

For each element of an input sequence X_i project into 3 vectors: query, key and value

For each element, compute attention over all other vectors

$$SA(Q_i, \mathbf{K}, \mathbf{V}) = \sum_j softmax_j \left(\frac{Q_i \cdot \mathbf{K}^T}{\sqrt{d_k}}\right) V_j$$

Vaswani et al., Attention Is All You Need, NIPS 2017

Structured Output and Attention Extending Memory Neural Reasoning

Self Attention – K,V,Q Generation

Self-attention

Figure credit to this article

Structured Output and Attention Extending Memory Neural Reasoning

Self Attention – Compute Attention Score

Self-attention

Structured Output and Attention Extending Memory Neural Reasoning

Self Attention – Produce Output

Self-attention

Structured Output and Attention Extending Memory Neural Reasoning

Self Attention – MultiHead

Strubell et al, Linguistically-Informed Self-Attention for Semantic Role Labeling, EMNLP 2018

Structured Output and Attention Extending Memory Neural Reasoning

Attention-Based Captioning – Focus Shifting

Soft Attention

Hard Attention

Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Structured Output and Attention Extending Memory Neural Reasoning

Attention-Based Captioning - Generation

Learns to correlate textual and visual concepts

A woman is throwing a <u>frisbee</u> in a park.

A dog is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

Helps understanding why the model fails

A large white bird standing in a forest.

A woman holding a <u>clock</u> in her hand.

Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Structured Output and Attention Extending Memory Neural Reasoning

Attention-Based Captioning – The Model

Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

RNN and Memory – Issue 1

- Gated RNN claim to solve the problem of learning long-range dependencies
- In practice it is still difficult to learn on longer range
- Architectures trying to optimize dynamic memory usage
 - Clockwork RNN
 - Skip RNN
 - Multiscale RNN
 - Zoneout

Clockwork RNN

Modular recurrent layer where each module is updated at different clock

Modules interconnected only when destination clock time is larger

Koutnik et al, A Clockwork RNN, ICML 2014

Skip RNN

fading effect due to the multiplicative gate)

Structured Output and Attention Extending Memory Neural Reasoning

Skip RNN and Attention

Attended pixels

Ignored pixels

Campos et al, Skip RNN: Skipping State Updates in Recurrent Neural Networks, ICLR 2018

Structured Output and Attention Extending Memory Neural Reasoning

RNN and Memory – Issue 2

A motivating example:

Task 3: Three Supporting Facts	Task 15: Basic Deduction
John picked up the apple.	Sheep are afraid of wolves.
John went to the office.	Cats are afraid of dogs.
John went to the kitchen.	Mice are afraid of cats.
John dropped the apple.	Gertrude is a sheep.
Where was the apple before the kitchen? A:office	What is Gertrude afraid of? A:wolves

- In order to solve the task need to memorize
 - Facts
 - Question
 - Answers
- A bit too much for the dynamical RNN memory
- Try to address it through an external memory

Memory Network Components

(I) **Input feature map**: Encodes the input in a feature vector

(G) **Generalization**: decide what input (or function of it) to write to memory

(O) **Output feature map**: reads the relevant memory slots

(R) **Response**: returns the prediction given the retrieved memories

IntroductionStructured Output and AttentionAdvanced RNNMemory NetworksWrap-upNeural Reasoning

End-to-End Memory Networks

Structured Output and Attention Extending Memory Neural Reasoning

Memory Network Extensions

Use more complex output components, e.g. RNN to generate response sequences

Stack multiple memory network layers

Several iterations of reasoning to produce a better answer

Memory Nets for Visual QA with Attention

Yang et al, Stacked Attention Networks for Image Question Answering, CVPR 2016

Structured Output and Attention Extending Memory Neural Reasoning

Memory Nets for Visual QA with Attention

(a) What are pulling a man on a wagon down on dirt road? Answer: horses Prediction: horses (b)

What is the color of the box ? Answer: red Prediction: red

What next to the large umbrella attached to a table? Answer: trees Prediction: tree

(c)

(d) How many people are going up the mountain with walking sticks? Answer: four Prediction: four

(e) What is sitting on the handle bar of a bicycle? Answer: bird Prediction: bird

What is the color of the horns? Answer: red Prediction: red

(f)

Structured Output and Attention Memory Networks Neural Reasoning

Neural Turing Machines

- Memory networks that can read and write memories at both training and test
- End-to-end differentiable

IntroductionStructured Output and AttentionAdvanced RNNMemory NetworksWrap-upNeural Reasoning

Neural Controller

Image credits @ colah.github.io

Typically an RNN emitting vectors to control read and write from the memory

The key to differentiability is to always read and write the whole memory

Use soft-attention to determine how much to read/write from each point IntroductionStructured Output and AttentionAdvanced RNNMemory NetworksWrap-upNeural Reasoning

Memory Read

Introduction Structu Advanced RNN Memor Wrap-up Neural

Structured Output and Attention Memory Networks Neural Reasoning

Memory Write

Structured Output and Attention Memory Networks Neural Reasoning

NTM Attention Focusing

1. Generate content-based memory indexing

Structured Output and Attention Memory Networks Neural Reasoning

NTM Attention Focusing

2. Interpolate with attention from previous time

Previous attention vector

Structured Output and Attention Memory Networks Neural Reasoning

NTM Attention Focusing

3. Generate location-based indexing

IntroductionStructured Output and AttentionAdvanced RNNMemory NetworksWrap-upNeural Reasoning

Practical Use?

- Not yet..
- Not straightforward to train
- Advantages over GRNN when it comes to learn to program
 - Copy task
 - Repeat copy
 - Associative recall
 - Sorting
- Foundations for neural reasoning
 - Pondering networks

Introduction Deep Gated RNN Wrap-up Software Conclusions

Software

- Complete sequence-to-sequence tutorial (including attention) on <u>Tensorflow</u>
 - A shorter version in <u>Keras</u>
- <u>Github project</u> collecting several memory augmented networks
- <u>Pytorch</u> implementation of stacked attention for visual QA (originally Theano-based)
- Many implementations of the NTM (Keras, Pytorch, Lasagne,...): none seemingly official, but <u>this recent one</u> is supposedly stable (enough for TF to list it as official)

Introduction Deep Gated RNN Wrap-up Software Conclusions

Take Home Messages

- Attention.. Attention.. and, again, attention
 - Soft attention is nice because makes everything fully differentiable
 - Hard attention is stochastic hence cannot Backprop
 - Empirical evidences of them being sensitive to different things
- Encoder-Decoder scheme
 - A general architecture to compose heterogeneous models and data
 - Decoding allows sampling complex predictions from an encoding conditioned distribution
- Memory and RNN
 - Efficient use of dynamic memory
 - External memory for search and recall tasks
 - Read/write memory for neural reasoning