1

Encoder-Decoder Sequence-to-Sequence Architectures

Sargur Srihari srihari@buffalo.edu

This is part of lecture slides on <a>Deep Learning: http://www.cedar.buffalo.edu/~srihari/CSE676

Topics in Sequence Modeling

- Overview
- 1. Unfolding Computational Graphs
- 2. Recurrent Neural Networks
- 3. Bidirectional RNNs
- 4. Encoder-Decoder Sequence-to-Sequence Architectures
- 5. Deep Recurrent Networks
- 6. Recursive Neural Networks
- 7. The Challenge of Long-Term Dependencies
- 8. Echo-State Networks
- 9. Leaky Units and Other Strategies for Multiple Time Scales
- 10. LSTM and Other Gated RNNs
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory



3

- 1. An RNN to map an input sequence to a fixed-size vector
- 2. RNN to map a fixed-size vector to a sequence
- 3. How an RNN can map an input sequence to an output sequence of the same length



RNN when input/output are not same length

- Here we consider how an RNN can be trained to map an input sequence to an output sequence which is not necessarily the same length
- Comes up in applications such as speech recognition, machine translation or question-answering where the input and output sequences in the training set are generally not of the same length (although their lengths might be related)

An Encoder-Decoder or Sequence-to-Sequence RNN



Learns to generate an output sequence $(\boldsymbol{y}^{(1)},..,\boldsymbol{y}^{(n_y)})$

given an input sequence

 $(m{x}^{(1)},..,m{x}^{(n_x)})$

It consists of an encoder RNN that reads an input sequence and a decoder RNN that generates the output sequence or computes the probability of a given output sequence)

The final hidden state of the encoder RNN Is used to compute a fixed size context *C* which represents a semantic summary of the input sequence and is given as jinput to the decoder