

Modeling sequences with recurrent networks

Julien VELCIN

Université Lumière Lyon 2

Master 2 MALIA

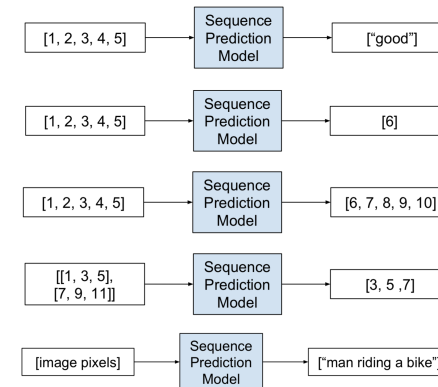
1 / 17

Some history

- aims to model more appropriately sequence structure
- foundational papers in both cognitive science and computational neuroscience journals:
 - Hopfield, 1982
 - Jordan, 1987
 - Elmann, 1990
- new era
 - LSTM
 - GRU
 - Bidirectional RNN

3 / 17

Multiple tasks



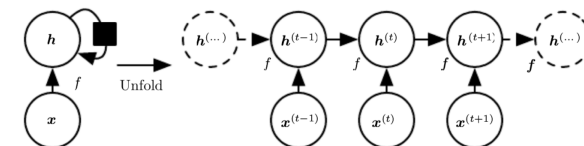
2 / 17

Simple RNN

Forward pass of a simple RNN at time t :

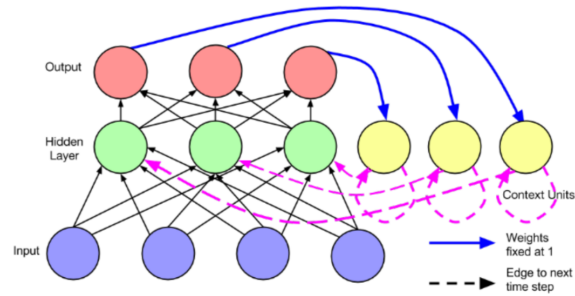
$$\mathbf{h}^{(t)} = f(W^{hx}\mathbf{x}^{(t)} + W^{hh}\mathbf{h}^{(t-1)} + \mathbf{b})$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W^{yh}\mathbf{h}^{(t)} + \mathbf{c})$$



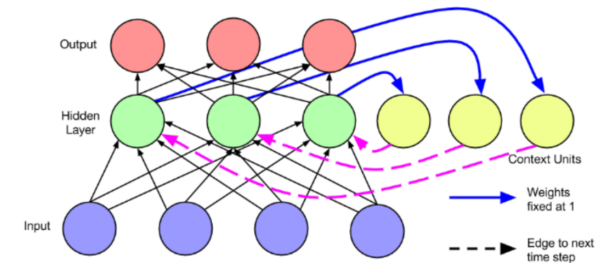
4 / 17

Jordan's net (1987)



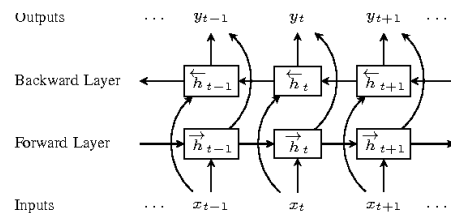
5 / 17

Elan's net (1990)



6 / 17

Bidirectional RNN (Schuster and Paliwal, 1997)



Picture taken from "A Unified Tagging Solution: Bidirectional LSTM Recurrent Neural Network with Word Embedding" (Wang et al., ArXiv 2015)

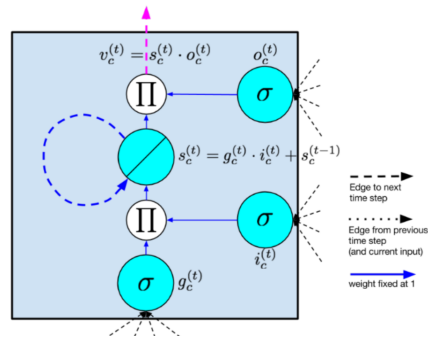
7 / 17

Problems of (past versions of) RNN

- training is difficult because optimisation is NP-complete
- long-range dependencies induces vanishing or exploding gradients
- truncated backprop through time can be one solution
 - cuts the time span influence
 - but kill the long-term memory
- a reborn of RNN was the introduction of gated architectures

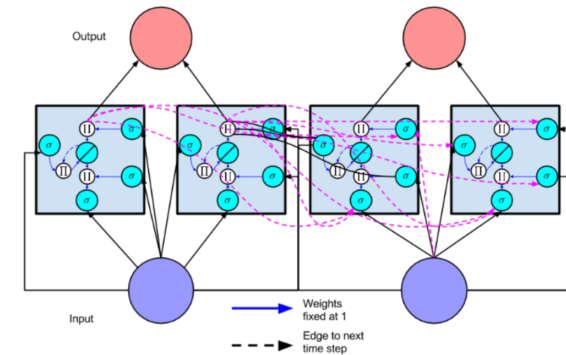
8 / 17

Long Short-Term Memory (Hochreiter and Schmidhuber, 1997)



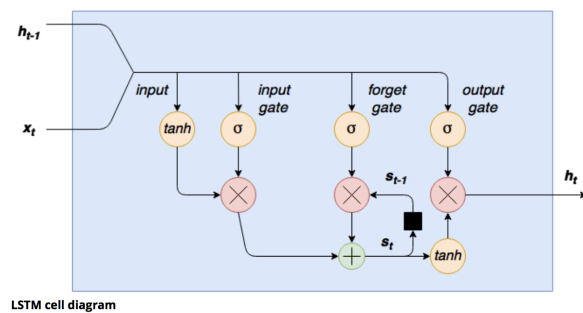
9 / 17

Long Short-Term Memory (unfolded)



10 / 17

Long Short-Term Memory (another view)

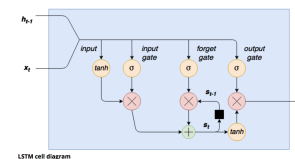


LSTM cell diagram

Taken from: <http://adventuresinmachinelearning.com/keras-lstm-tutorial/>

11 / 17

LSTM (math)



$$g = \tanh(b^g + x_t \cdot U^g + h_{t-1} \cdot V^g)$$

$$i = \sigma(b^i + x_t \cdot U^i + h_{t-1} \cdot V^i)$$

$$f = \sigma(b^f + x_t \cdot U^f + h_{t-1} \cdot V^f)$$

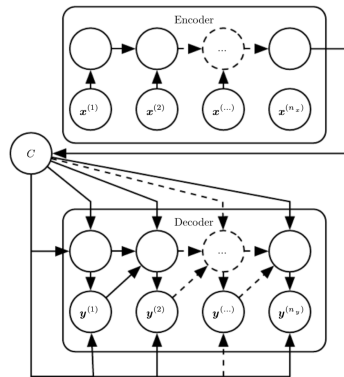
$$s_t = s_{t-1} \odot f + g \odot i$$

$$o = \sigma(b^o + x_t \cdot U^o + h_{t-1} \cdot V^o)$$

$$h_t = \tanh(s_t) \odot o$$

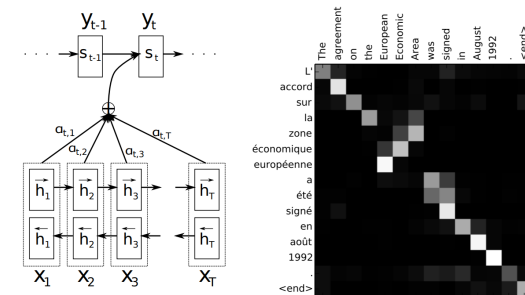
12 / 17

seq2seq architectures (Cho et al., 2014; Sutskever et al., 2014)



13 / 17

All you need is attention? (Bahdanau et al., 2015)



14 / 17

The Transformer (Vaswani et al., 2017)

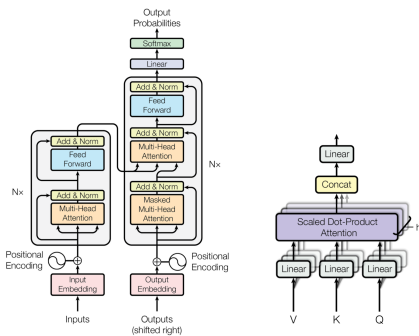
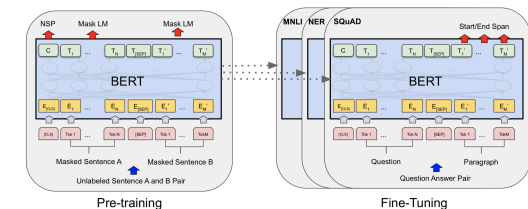
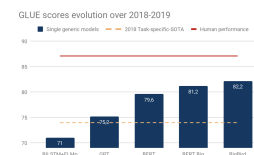


Figure 1: The Transformer - model architecture.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

15 / 17

BERT and why attention isn't all you need



Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of NAACL-HLT 2019*, pages 4171-4186.

16 / 17

Some references

- Elman, Jeffrey L. (1990). Finding Structure in Time. Cognitive Science. 14 (2): 179–211.
- Hochreiter, S. and Schmidhuber, J. (1997). LSTM can solve hard long time lag problems. In Advances in neural information processing systems (pp. 473-479).
- Bahdanau, D., Cho, K. and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. ICLR 2015.
- Goodfellow, I., Bengio, Y., Courville, A. and Bengio, Y. (2016). Deep learning. Cambridge: MIT press.
- Brownlee, J. (2017). Long Short-Term Memory Networks with Python Develop Sequence Prediction Models with Deep Learning. Machine Learning Mastery, EBook.