

Chapter 10

Sequence Modelling: Recurrent and Recursive Nets

10.1 Unfolding Computational Graphs

10.1.1

The authors state that "Much as almost any function can be considered a feed-forward neural network, essentially any function involving recurrence can be considered a recurrent neural network." Do you agree with this statement?

10.1.2

When predicting the future from the past, an RNN learns to use its hidden state to summarize the history of observations. This is in general a necessarily lossy summary, since the history sequence can be any length and the hidden state is finite. Give an example when this is irrelevant for the prediction task. Give an example when this makes a huge difference.

10.1.3

The unfolded view of a neural network illustrates how information flows forward in time. Is this necessarily the same time direction as in the data (i.e. from an observed time-series)? If yes, explain why. If no, give a counter-example.

10.1.4

If we were to imagine neural networks as implementing algorithms using a Control-Flow Graph, how would you characterize the difference between the CFG implemented by a Feedforward Neural Network compared to a Recurrent Neural Network?

10.2 Recurrent Neural Networks

10.2.1

The authors assume hyperbolic tangent units, not ReLU which in earlier chapters has been the most common one. Why do you think that is?

10.2.2

What is the fundamental difference between the two deep RNNs described by the equations 10.1 and 10.2 below?

$$\begin{aligned}
 \mathbf{h}_1^{(t)} &= \tanh(\mathbf{b}_1 + W_1 \mathbf{h}_1^{(t-1)} + U_1 \mathbf{x}^{(t)}) \\
 \mathbf{h}_2^{(t)} &= \tanh(\mathbf{b}_2 + W_2 \mathbf{h}_2^{(t-1)} + U_2 \mathbf{h}_1^{(t)}) \\
 &\vdots \\
 \mathbf{h}_l^{(t)} &= \tanh(\mathbf{b}_l + W_l \mathbf{h}_l^{(t-1)} + U_l \mathbf{h}_{l-1}^{(t)}) \\
 \mathbf{o}^{(t)} &= c + V \mathbf{h}_l^{(t)}
 \end{aligned} \tag{10.1}$$

$$\begin{aligned}
 \mathbf{h}_1^{(t)} &= \tanh(\mathbf{b}_1 + W_1 \mathbf{x}^{(t-1)} + U_1 \mathbf{x}^{(t)}) \\
 \mathbf{h}_2^{(t)} &= \tanh(\mathbf{b}_2 + W_2 \mathbf{h}_1^{(t-1)} + U_2 \mathbf{h}_1^{(t)}) \\
 &\vdots \\
 \mathbf{h}_l^{(t)} &= \tanh(\mathbf{b}_l + W_l \mathbf{h}_{l-1}^{(t-1)} + U_l \mathbf{h}_{l-1}^{(t)}) \\
 \mathbf{o}^{(t)} &= c + V \mathbf{h}_l^{(t)}
 \end{aligned} \tag{10.2}$$

10.2.3

Below (10.3) is a reproduction of equation 10.14 from the book. Under what independence assumptions is it modelling the joint probability $P(y^{(1)}, y^{(2)}, \dots, y^{(1)})$?

$$- \sum_t \log p_{\text{model}}(y^{(t)} | \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}\}) \tag{10.3}$$

10.2.4

A common setup for training RNNs on sequence data (text and time-series) is to use next-step prediction (see equation 10.4 below). In what way (if any) is this different from the network described by figure 10.4 of the book, and is it an application of teacher forcing?

$$\begin{aligned}
 \mathbf{h}^{(t)} &= \tanh(\mathbf{b} + W \mathbf{h}^{(t-1)} + U \mathbf{x}^{(t)}) \\
 \mathbf{o}^{(t)} &= c + V \mathbf{h}^{(t)} \\
 \mathbf{y}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}) \\
 L^{(t)} &= -\log p_{\text{model}}(\mathbf{x}^{(t+1)} | \mathbf{o}^{(t)})
 \end{aligned} \tag{10.4}$$