

time steps. Without this extra input, the RNN might generate sequences that end abruptly, such as a sentence that ends before it is complete. This approach is based on the decomposition

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}) = P(\tau)P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)} \mid \tau). \quad (10.34)$$

The strategy of predicting τ directly is used for example by Goodfellow *et al.* (2014d).

10.2.4 Modeling Sequences Conditioned on Context with RNNs

In the previous section we described how an RNN could correspond to a directed graphical model over a sequence of random variables $y^{(t)}$ with no inputs \mathbf{x} . Of course, our development of RNNs as in equation 10.8 included a sequence of inputs $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(\tau)}$. In general, RNNs allow the extension of the graphical model view to represent not only a joint distribution over the y variables but also a conditional distribution over y given \mathbf{x} . As discussed in the context of feedforward networks in section 6.2.1.1, any model representing a variable $P(\mathbf{y}; \boldsymbol{\theta})$ can be reinterpreted as a model representing a conditional distribution $P(\mathbf{y}|\boldsymbol{\omega})$ with $\boldsymbol{\omega} = \boldsymbol{\theta}$. We can extend such a model to represent a distribution $P(\mathbf{y} \mid \mathbf{x})$ by using the same $P(\mathbf{y} \mid \boldsymbol{\omega})$ as before, but making $\boldsymbol{\omega}$ a function of \mathbf{x} . In the case of an RNN, this can be achieved in different ways. We review here the most common and obvious choices.

Previously, we have discussed RNNs that take a sequence of vectors $\mathbf{x}^{(t)}$ for $t = 1, \dots, \tau$ as input. Another option is to take only a single vector \mathbf{x} as input. When \mathbf{x} is a fixed-size vector, we can simply make it an extra input of the RNN that generates the \mathbf{y} sequence. Some common ways of providing an extra input to an RNN are:

1. as an extra input at each time step, or
2. as the initial state $\mathbf{h}^{(0)}$, or
3. both.

The first and most common approach is illustrated in figure 10.9. The interaction between the input \mathbf{x} and each hidden unit vector $\mathbf{h}^{(t)}$ is parametrized by a newly introduced weight matrix \mathbf{R} that was absent from the model of only the sequence of y values. The same product $\mathbf{x}^\top \mathbf{R}$ is added as additional input to the hidden units at every time step. We can think of the choice of \mathbf{x} as determining the value