

## Training RNNs

To train an RNN, the trick is to unroll it through time (like we just did) and then simply use regular backpropagation (see [Figure 14-5](#)). This strategy is called *backpropagation through time* (BPTT).

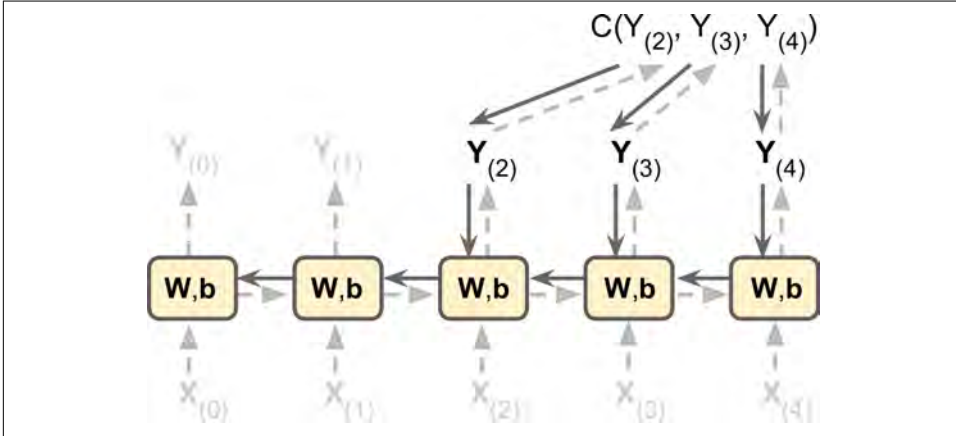


Figure 14-5. Backpropagation through time

Just like in regular backpropagation, there is a first forward pass through the unrolled network (represented by the dashed arrows); then the output sequence is evaluated using a cost function  $C(Y_{(t_{\min})}, Y_{(t_{\min} + 1)}, \dots, Y_{(t_{\max})})$  (where  $t_{\min}$  and  $t_{\max}$  are the first and last output time steps, not counting the ignored outputs), and the gradients of that cost function are propagated backward through the unrolled network (represented by the solid arrows); and finally the model parameters are updated using the gradients computed during BPTT. Note that the gradients flow backward through all the outputs used by the cost function, not just through the final output (for example, in [Figure 14-5](#) the cost function is computed using the last three outputs of the network,  $Y_{(2)}$ ,  $Y_{(3)}$ , and  $Y_{(4)}$ , so gradients flow through these three outputs, but not through  $Y_{(0)}$  and  $Y_{(1)}$ ). Moreover, since the same parameters  $W$  and  $b$  are used at each time step, backpropagation will do the right thing and sum over all time steps.

## Training a Sequence Classifier

Let's train an RNN to classify MNIST images. A convolutional neural network would be better suited for image classification (see [Chapter 13](#)), but this makes for a simple example that you are already familiar with. We will treat each image as a sequence of 28 rows of 28 pixels each (since each MNIST image is  $28 \times 28$  pixels). We will use cells of 150 recurrent neurons, plus a fully connected layer containing 10 neurons