

Recurrent neurons (p.381)

phunc20

January 28, 2021

Abstract

Personally, I found the equations on p.381 were not to the point, at least not enough for me to grasp the exact meaning and shapes of the input and output of a layer of recurrent neurons. Every time I re-read these paragraphs, I cannot recall the understanding I reached at the previous reading. Thus, I decided to make a note here to try to change the situation, hoping to make things clearer.

1 A Single Recurrent Neuron

The output of a single recurrent neuron of a single input instance $x_{(t)}$ equals

$$y_{(t)} = \phi\left(\langle x_{(t)}, w_x \rangle + y_{(t-1)}w_y + b\right),$$

where

- $x_{(t)}, w_x \in \mathbb{R}^k$ (k can be any positive integer the user consider suitable. Note that k is called $\#(\text{inputs})$ in the textbook)
- $\langle \cdot, \cdot \rangle$ denotes the inner product in \mathbb{R}^k
- $y_{(t-1)} \in \mathbb{R}$ is the output of the previous time step
- w_y, b are both in \mathbb{R}

Recurrent Neurons

Up to now we have mostly looked at feedforward neural networks, where the activations flow only in one direction, from the input layer to the output layer (except for a few networks in [Appendix E](#)). A recurrent neural network looks very much like a feedforward neural network, except it also has connections pointing backward. Let's look at the simplest possible RNN, composed of just one neuron receiving inputs, producing an output, and sending that output back to itself, as shown in [Figure 14-1](#) (left). At each *time step* t (also called a *frame*), this *recurrent neuron* receives the inputs $\mathbf{x}_{(t)}$ as well as its own output from the previous time step, $y_{(t-1)}$. We can represent this tiny network against the time axis, as shown in [Figure 14-1](#) (right). This is called *unrolling the network through time*.

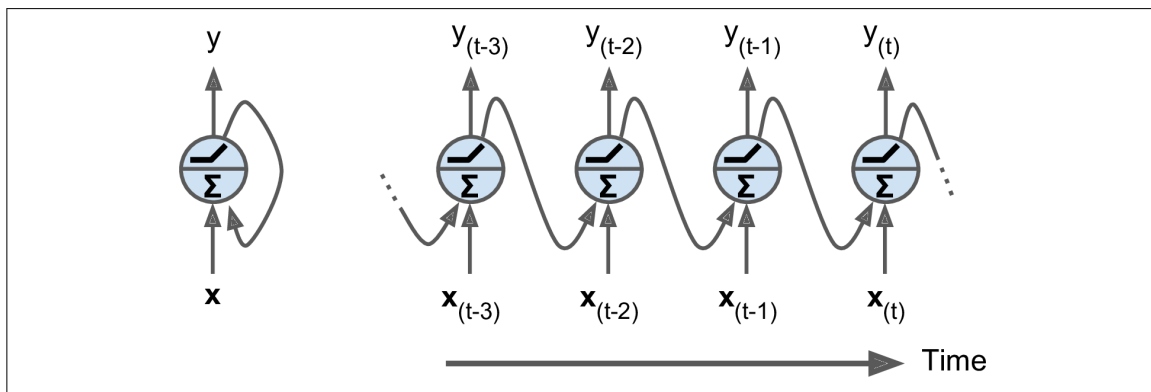


Figure 14-1. A recurrent neuron (left), unrolled through time (right)

You can easily create a layer of recurrent neurons. At each time step t , every neuron receives both the input vector $\mathbf{x}_{(t)}$ and the output vector from the previous time step $\mathbf{y}_{(t-1)}$, as shown in [Figure 14-2](#). Note that both the inputs and outputs are vectors now (when there was just a single neuron, the output was a scalar).

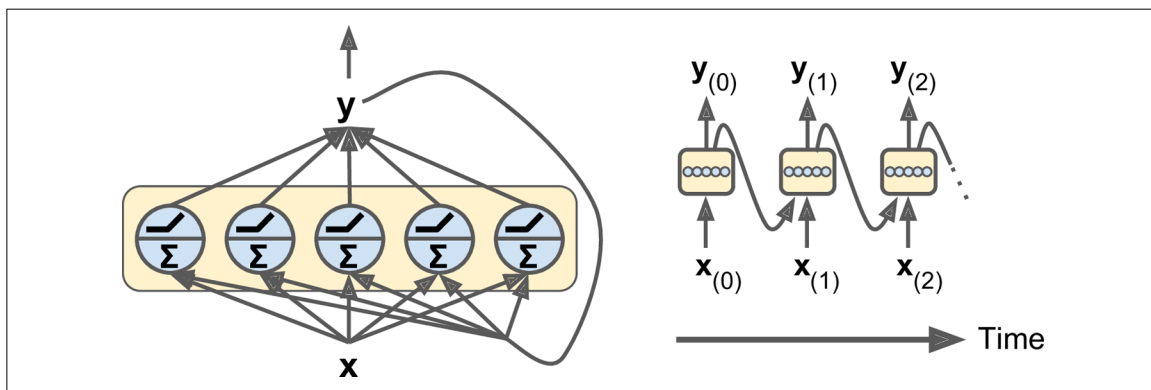


Figure 14-2. A layer of recurrent neurons (left), unrolled through time (right)

Each recurrent neuron has two sets of weights: one for the inputs $\mathbf{x}_{(t)}$ and the other for the outputs of the previous time step, $\mathbf{y}_{(t-1)}$. Let's call these weight vectors \mathbf{w}_x and \mathbf{w}_y .

The output of a single recurrent neuron can be computed pretty much as you might expect, as shown in [Equation 14-1](#) (b is the bias term and $\phi(\cdot)$ is the activation function, e.g., ReLU¹).

Equation 14-1. Output of a single recurrent neuron for a single instance

$$\mathbf{y}_{(t)} = \phi\left(\mathbf{x}_{(t)}^T \cdot \mathbf{w}_x + \mathbf{y}_{(t-1)}^T \cdot \mathbf{w}_y + b\right)$$

Just like for feedforward neural networks, we can compute a whole layer's output in one shot for a whole mini-batch using a vectorized form of the previous equation (see [Equation 14-2](#)).

Equation 14-2. Outputs of a layer of recurrent neurons for all instances in a mini-batch

$$\begin{aligned}\mathbf{Y}_{(t)} &= \phi\left(\mathbf{X}_{(t)} \cdot \mathbf{W}_x + \mathbf{Y}_{(t-1)} \cdot \mathbf{W}_y + \mathbf{b}\right) \\ &= \phi\left(\begin{bmatrix} \mathbf{X}_{(t)} & \mathbf{Y}_{(t-1)} \end{bmatrix} \cdot \mathbf{W} + \mathbf{b}\right) \text{ with } \mathbf{W} = \begin{bmatrix} \mathbf{W}_x \\ \mathbf{W}_y \end{bmatrix}\end{aligned}$$

- $\mathbf{Y}_{(t)}$ is an $m \times n_{\text{neurons}}$ matrix containing the layer's outputs at time step t for each instance in the mini-batch (m is the number of instances in the mini-batch and n_{neurons} is the number of neurons).
- $\mathbf{X}_{(t)}$ is an $m \times n_{\text{inputs}}$ matrix containing the inputs for all instances (n_{inputs} is the number of input features).
- \mathbf{W}_x is an $n_{\text{inputs}} \times n_{\text{neurons}}$ matrix containing the connection weights for the inputs of the current time step.
- \mathbf{W}_y is an $n_{\text{neurons}} \times n_{\text{neurons}}$ matrix containing the connection weights for the outputs of the previous time step.
- The weight matrices \mathbf{W}_x and \mathbf{W}_y are often concatenated into a single weight matrix \mathbf{W} of shape $(n_{\text{inputs}} + n_{\text{neurons}}) \times n_{\text{neurons}}$ (see the second line of [Equation 14-2](#)).
- \mathbf{b} is a vector of size n_{neurons} containing each neuron's bias term.

¹ Note that many researchers prefer to use the hyperbolic tangent (tanh) activation function in RNNs rather than the ReLU activation function. For example, take a look at by Vu Pham et al.'s paper "[Dropout Improves Recurrent Neural Networks for Handwriting Recognition](#)". However, ReLU-based RNNs are also possible, as shown in Quoc V. Le et al.'s paper "[A Simple Way to Initialize Recurrent Networks of Rectified Linear Units](#)".

2 A Layer of Recurrent Neurons

To illustrate, let's take Figure 14-2 for example, in which we have a layer of **5** recurrent neurons. So, the output of this particular layer equals

$$\begin{bmatrix} y_{(t),1} \\ y_{(t),2} \\ y_{(t),3} \\ y_{(t),4} \\ y_{(t),5} \end{bmatrix} = \phi \left(\begin{bmatrix} \langle x_{(t)}, w_{x,1} \rangle + \langle y_{(t-1)}, w_{y,1} \rangle + b_1 \\ \langle x_{(t)}, w_{x,2} \rangle + \langle y_{(t-1)}, w_{y,2} \rangle + b_2 \\ \langle x_{(t)}, w_{x,3} \rangle + \langle y_{(t-1)}, w_{y,3} \rangle + b_3 \\ \langle x_{(t)}, w_{x,4} \rangle + \langle y_{(t-1)}, w_{y,4} \rangle + b_4 \\ \langle x_{(t)}, w_{x,5} \rangle + \langle y_{(t-1)}, w_{y,5} \rangle + b_5 \end{bmatrix} \right), \text{ i.e.}$$

$$y_{(t)} = \phi \left(W_x^T x_{(t)} + W_y^T y_{(t-1)} + b \right),$$

where

- $y_{(t),j} \in \mathbb{R}$ for all $j = 1, 2, 3, 4, 5$
- $x_{(t)} \in \mathbb{R}^k$ as before, and $y_{(t-1)} \in \mathbb{R}^l$, where $l = \#(\text{recurrent neurons})$, i.e. number of recurrent neurons in the layer, here, in particular, $l = 5$
- $w_{x,j} \in \mathbb{R}^k, w_{y,j} \in \mathbb{R}^l, b_j \in \mathbb{R}$ are the weights and bias associated with the j -th neuron (for all $j = 1, 2, 3, 4, 5$)
- We define

$$W_x^T = \begin{bmatrix} \text{--- } w_{x,1}^T \text{ ---} \\ \text{--- } w_{x,2}^T \text{ ---} \\ \text{--- } w_{x,3}^T \text{ ---} \\ \text{--- } w_{x,4}^T \text{ ---} \\ \text{--- } w_{x,5}^T \text{ ---} \end{bmatrix}$$

or equivalently

$$W_x = \begin{bmatrix} | & | & | & | & | \\ w_{x,1} & w_{x,2} & w_{x,3} & w_{x,4} & w_{x,5} \\ | & | & | & | & | \end{bmatrix}.$$

- Defining W_x and using the transpose this way has a similar reason as before when we were dealing with linear regression or fully connected NN
 - Because later on we will use batch and have matrix multiplication XW_x .

- We define similarly

$$W_y^T = \begin{bmatrix} \text{---} w_{y,1}^T \text{---} \\ \text{---} w_{y,2}^T \text{---} \\ \text{---} w_{y,3}^T \text{---} \\ \text{---} w_{y,4}^T \text{---} \\ \text{---} w_{y,5}^T \text{---} \end{bmatrix}$$

or equivalently

$$W_y = \begin{bmatrix} | & | & | & | & | \\ w_{y,1} & w_{y,2} & w_{y,3} & w_{y,4} & w_{y,5} \\ | & | & | & | & | \end{bmatrix}.$$

- It is understood that

$$b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{bmatrix}$$

3 A Layer of Recurrent Neurons on A Batch of Instances

Like in the case of a fully connected neural network, we can have an expression (using matrices) for an entire batch of input instances, say, m instances

$$\begin{aligned} Y_{(t)} &= \phi \left(X_{(t)} W_x + Y_{(t-1)} W_y + b \right) \\ &= \phi \left(\begin{bmatrix} X_{(t)} & Y_{(t-1)} \end{bmatrix} W + b \right) \text{ with } W = \begin{bmatrix} W_x \\ W_y \end{bmatrix} \end{aligned}$$

where

- $\text{shape}(X_{(t)}) = (m, k)$
- $\text{shape}(W_x) = (k, l)$
- $\text{shape}(Y_{(t-1)}) = (m, l)$
- $\text{shape}(W_y) = (l, l)$
- $\text{shape}(b) = (1, l)$