deeplearning.ai

Recurrent Neural Networks

Recurrent Neural Network Model

# Why not a standard network?

$$x^{<1>}$$
$$x^{<2>}$$
$$x^{<t:>}$$
$$x^{<T_x>}$$

$$y^{<1>}$$
$$y^{<2>}$$
$$y^{<T_y>}$$

pad / zero pad

This is maybe similar to what you saw in convolutional neural network where you want things learned for one part of the image to generalize quickly to the other parts of the image, and we'd like similar effect for sequence data as well. And similar to what you saw with convnets, using a better representation will also let you reduce the number of parameters in your model.
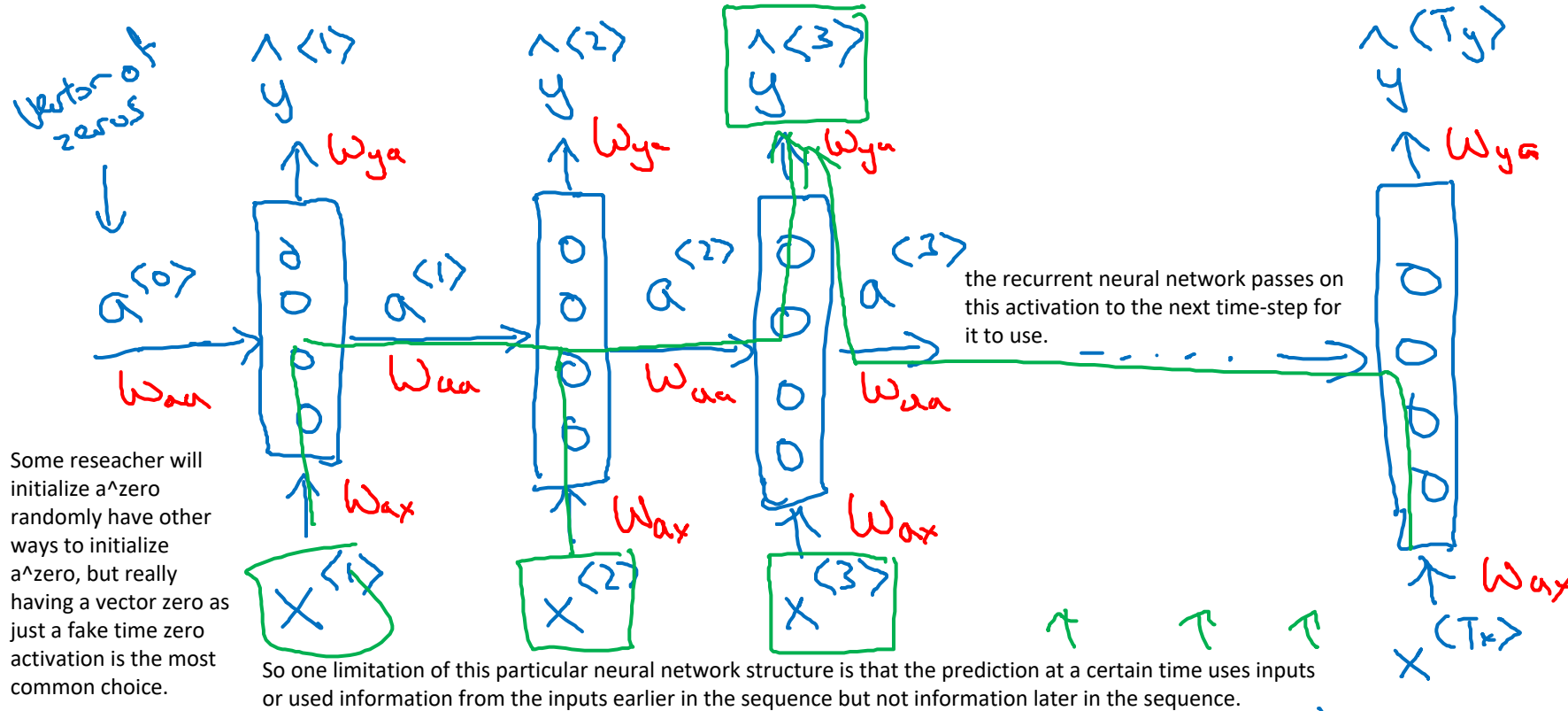
## Problems:

- Inputs, outputs can be different lengths in different examples.

- Doesn't share features learned across different positions of text.

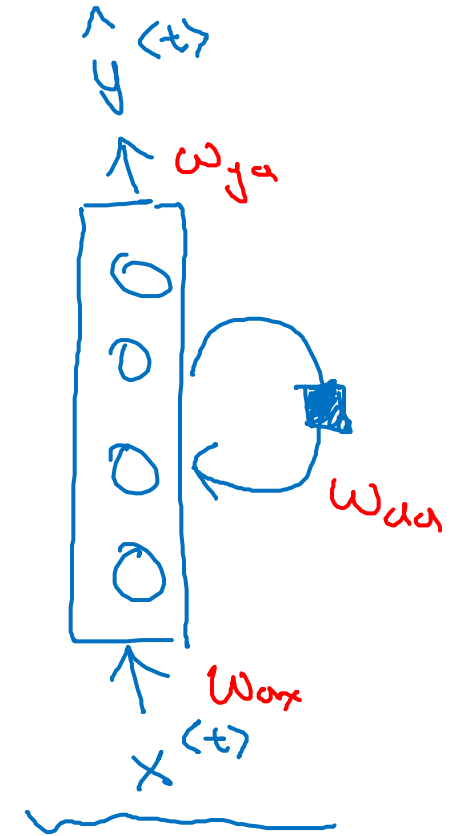Andrew Ng

# Recurrent Neural Networks

$T_x = T_y$

The architecture will change a little bit if T_x and T_y are not identical.

the recurrent neural network passes on this activation to the next time-step for it to use.

Some reseacher will initialize a^zero randomly have other ways to initialize a^zero, but really having a vector zero as just a fake time zero activation is the most common choice.

So one limitation of this particular neural network structure is that the prediction at a certain time uses inputs or used information from the inputs earlier in the sequence but not information later in the sequence.
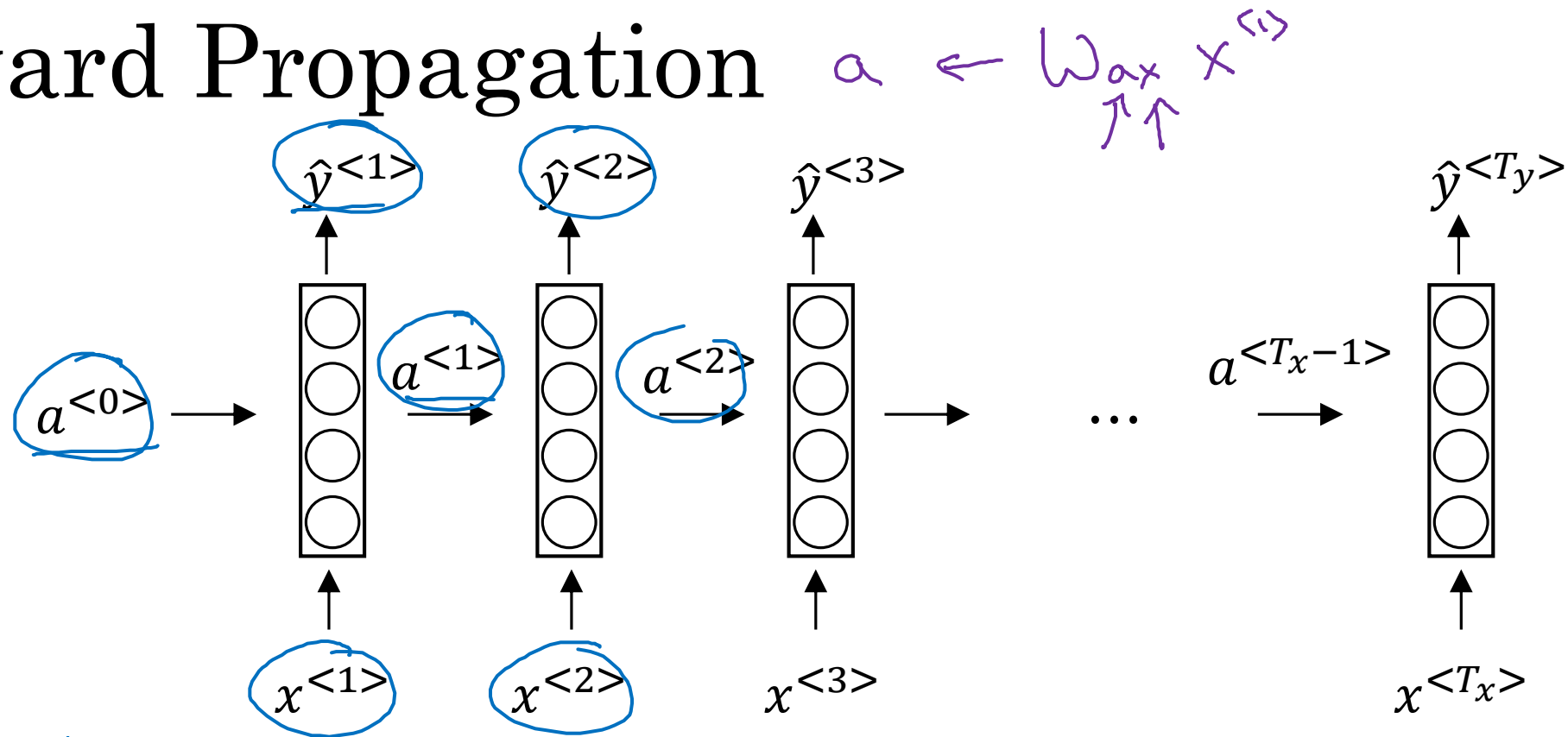
Bidirectional RNN (BRNN)

He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

Andrew Ng

# Forward Propagation

$$a \leftarrow W_{ax} \, x^{\langle 1 \rangle}$$



$$a^{\langle 0 \rangle} = \vec{0}.$$

$$a^{\langle 1 \rangle} = g_1 \left( W_{aa} \, a^{\langle 0 \rangle} + W_{ax} \, x^{\langle 1 \rangle} + b_a \right) \quad \leftarrow \text{tanh} \, / \, \text{ReLU}$$

tanh is actually a pretty common choice we have other ways of preventing the vanishing gradient problem.

$$\hat{y}^{\langle 1 \rangle} = g_2 \left( W_{ya} \, a^{\langle 1 \rangle} + b_y \right) \quad \leftarrow \text{Sigmoid}$$

$$a^{\langle t \rangle} = g \left( W_{aa} \, a^{\langle t-1 \rangle} + W_{ax} \, x^{\langle t \rangle} + b_a \right)$$

$$\hat{y}^{\langle t \rangle} = g \left( W_{ya} \, a^{\langle t \rangle} + b_y \right)$$

Andrew Ng

# Simplified RNN notation

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$(100,100)$   100   $(100,10,000)$   10,000

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

$$y^{<t>} = g(W_y a^{<t>} + b_y)$$

$$a^{<t>} = g\left(W_a [a^{<t-1>}, x^{<t>}] + b_a\right)$$

$$\begin{bmatrix} W_{aa} & ; & W_{ax} \end{bmatrix} = W_a$$
100   10000

$(100, 10100)$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \begin{matrix} 100 \\ 10000 \end{matrix} \; 10100$$

$$[W_{aa} ; W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$