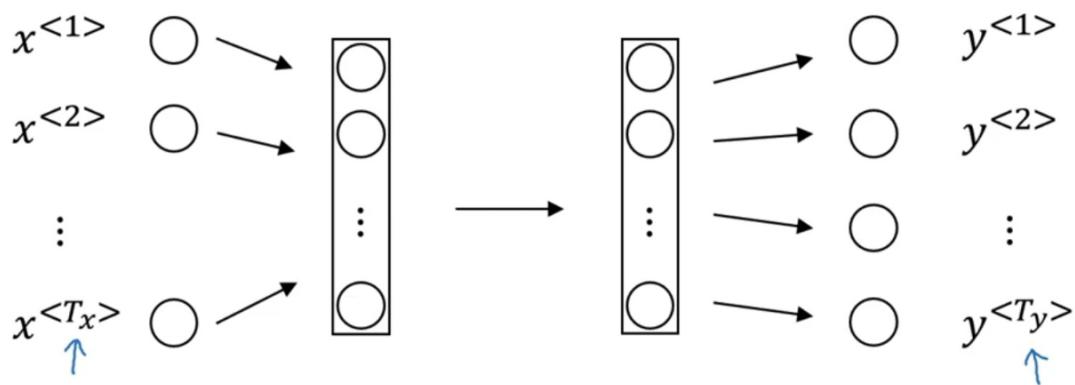


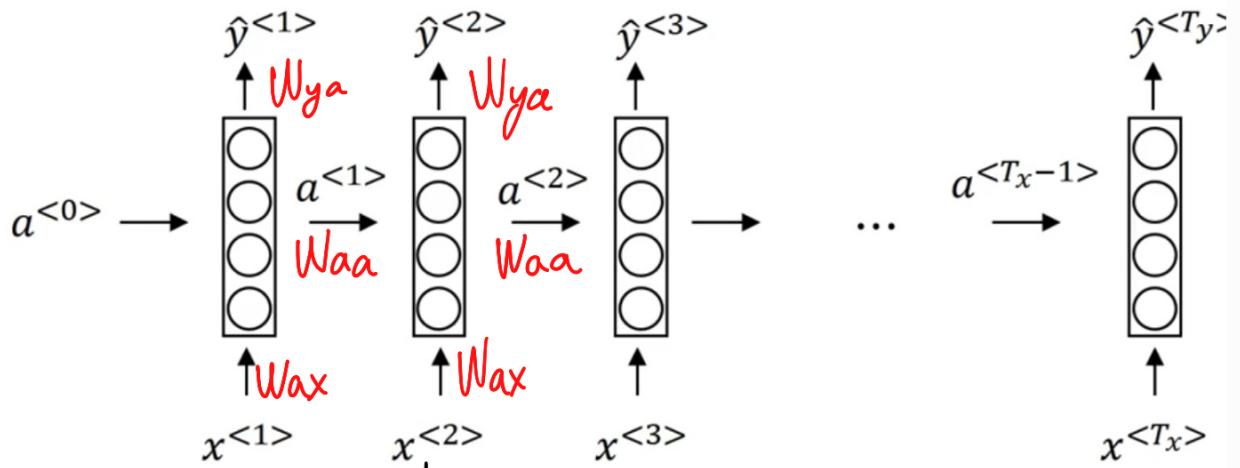
Why not a standard network?



Problems:

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Forward Propagation



A RNN, computes each $a^{<i>}$ and $\hat{y}^{<i>}$, using $x^{<i-1>} \rightarrow$ and $a^{<i-1>} \rightarrow$, which defines a recursive function with base case $a^{<0>} \rightarrow$ and $x^{<1>} \rightarrow$.

However, even though it is great that it uses past information, it does not use future one. Therefore, RNN would not be able to distinguish who is a person here:

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

He said, "Teddy bears are on sale!"

$$a^{<t>} = g(W_{aa}a^{<t>} + W_{ax}x^{<t>} + b_a) \quad \text{let's simplify:}$$

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

$$Wa = [W_{aa} : W_{ax}] \quad (\text{they both have same number of rows})$$

$\hookrightarrow \in \mathbb{R}^{n\text{-rows_both} \times (n\text{-cols_} Wa + n\text{-cols_} W_{ax})}$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ \dots \\ x^{<t>} \end{bmatrix} \in \mathbb{R}^{(n\text{-rows_} a^{<t-1>} + n\text{-rows_} x^{<t>}) \times 1}$$

So, we rewrite $a^{<t>}$ as:

$$a^{<t>} = g(Wa [a^{<t-1>}, x^{<t>}] + ba)$$

And $\hat{y}^{<t>}$ as: usually tanh, less usually ReLU

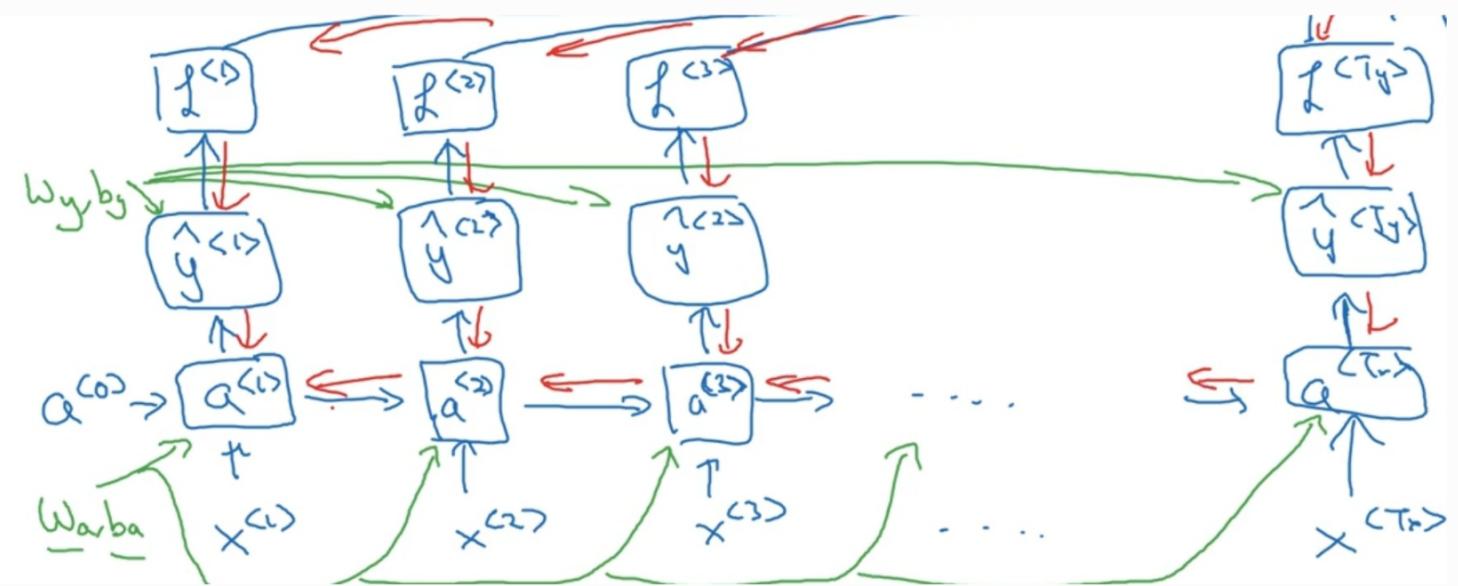
$$\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$

Loss function: softmax, sigmoid

$L^{<t>}(\hat{y}^{<t>}, y^{<t>})$ = cross-entropy loss, from logistic regression =

$$-y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

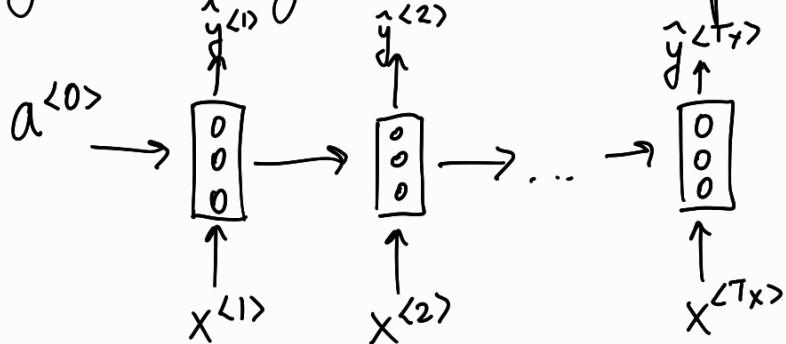
$$[L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})]$$



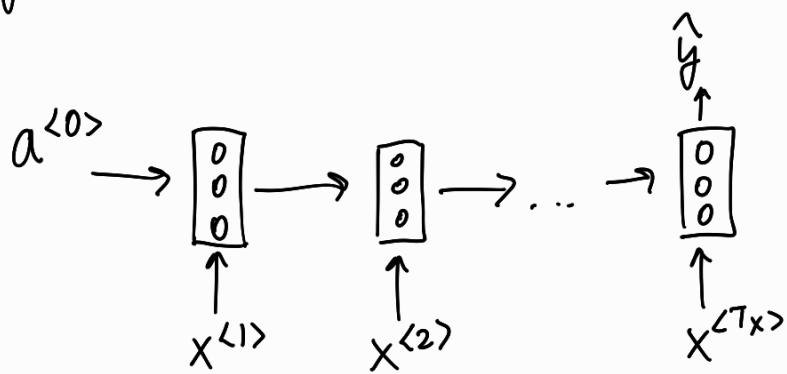
Backpropagation through time.

RNNs architectures:

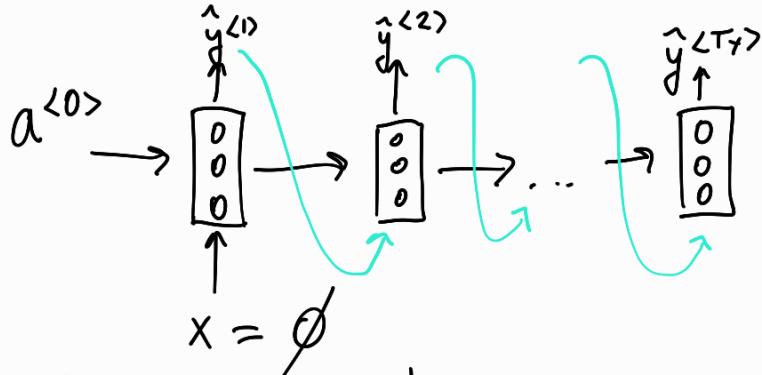
- Many-to-many (the one exposed above):



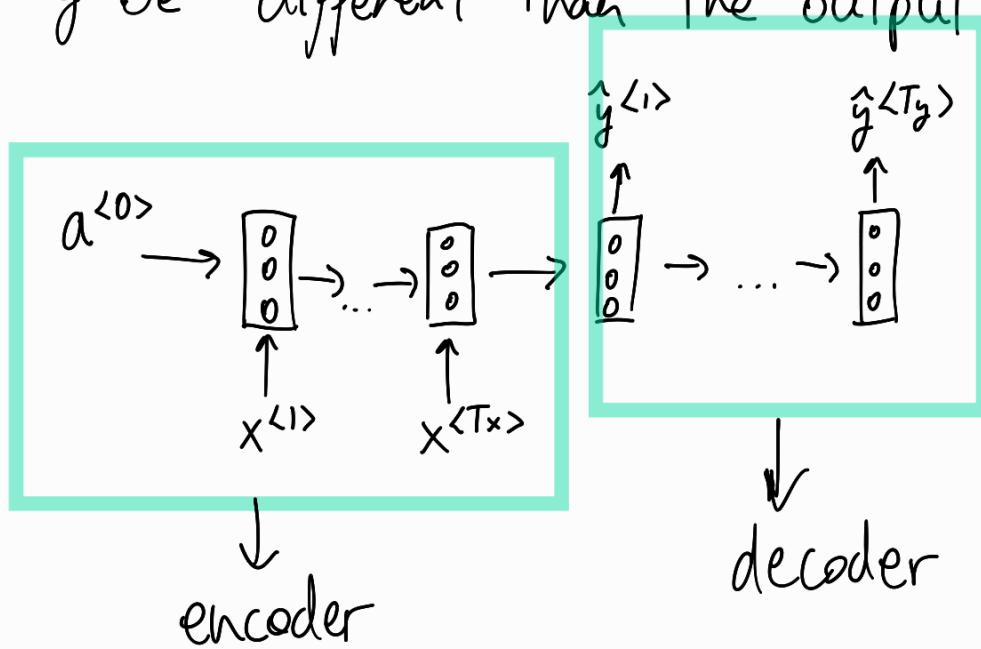
- Many-to-one: for sentiment classification for example



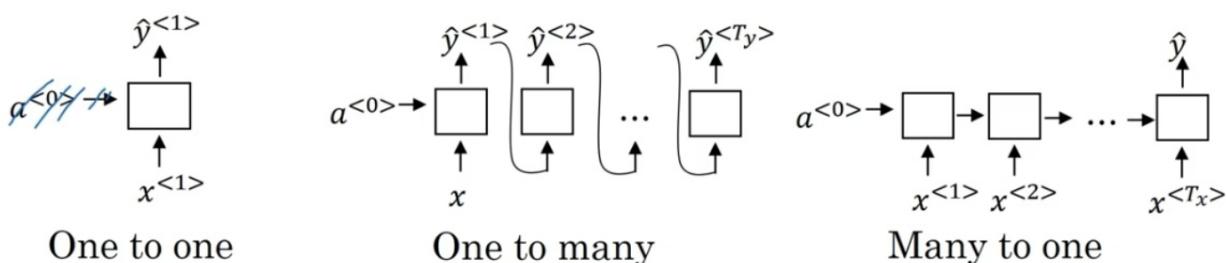
- One-to-many: for example, music generation.



- Many-to-many, with different lengths for input and output. For example, in machine translation, the number of words in the input sentence (in French) may be different than the output one (in English).



Summary of RNN types



Vanishing gradients and forgetting information:

Each $y^{<i>}$, is mainly influenced by close $x^{<j>}$ to it, $j < i$.

That means, the RNN "forgets" as it advances in sequences. So for instance, when generating sentences like:

"the cat, ..., was full"
"The cats, ..., were full"

The RNN, may forget about if there was 1 or more cats.

That is because of backpropagation.

- Solution : GRUs

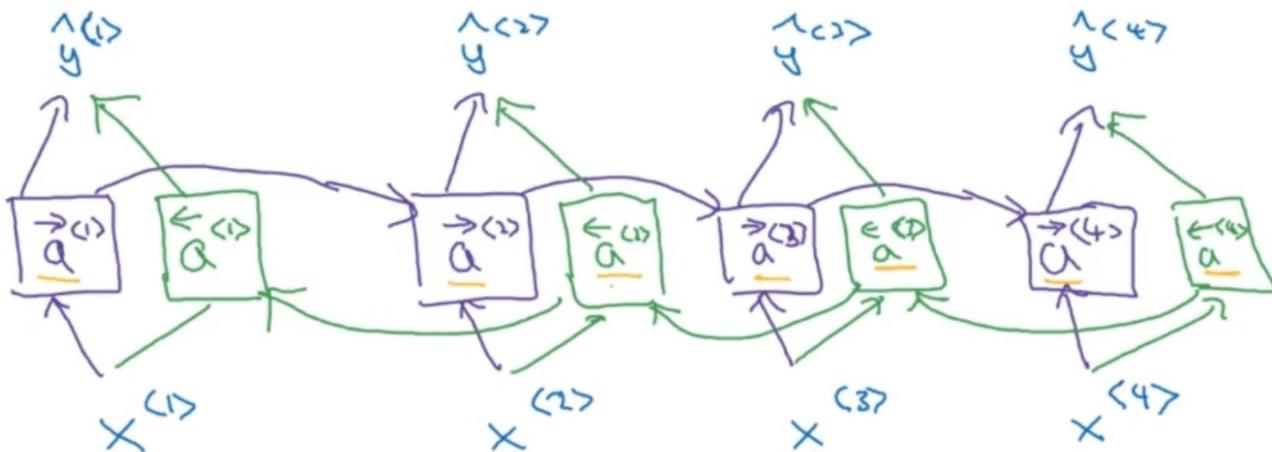
Exploding gradients: Use gradient clipping.

Bidirectional RNNs:

Getting information from the future

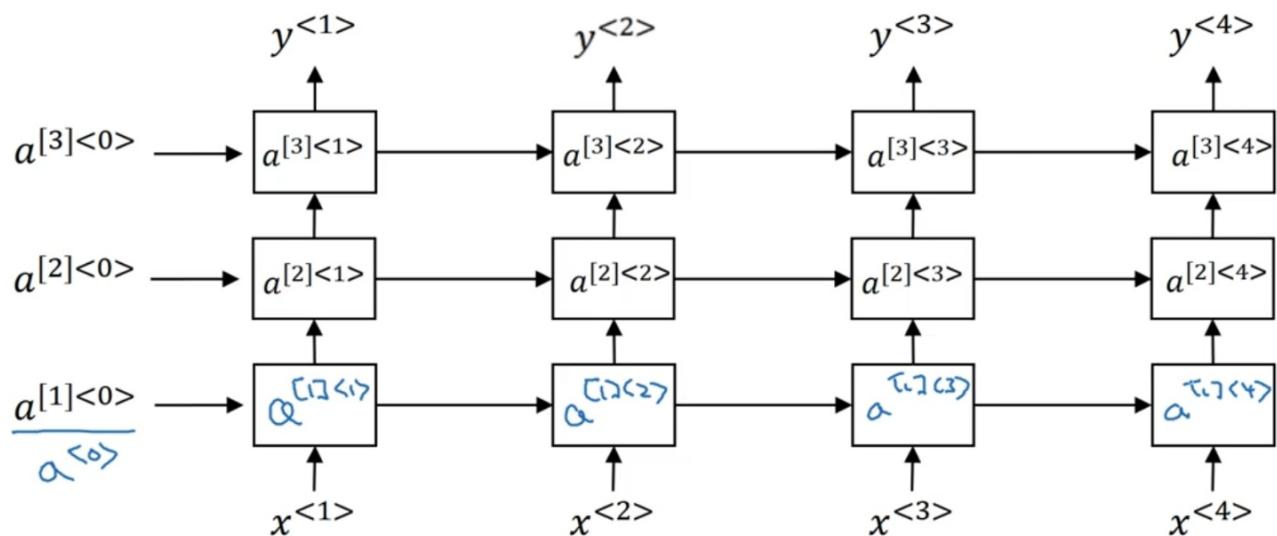
He said, "Teddy bears are on sale!"

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



Disadvantage : we need the entire sequence before we can predict.

Deep RNNs:



$$a^{[2]<3>} = g(W_a^{[2]} [a^{[1]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

What you should remember:

- The recurrent neural network, or RNN, is essentially the repeated use of a single cell.
- A basic RNN reads inputs one at a time, and remembers information through the hidden layer activations (hidden states) that are passed from one time step to the next.
 - The time step dimension determines how many times to re-use the RNN cell
- Each cell takes two inputs at each time step:
 - The hidden state from the previous cell
 - The current time step's input data
- Each cell has two outputs at each time step:
 - A hidden state
 - A prediction

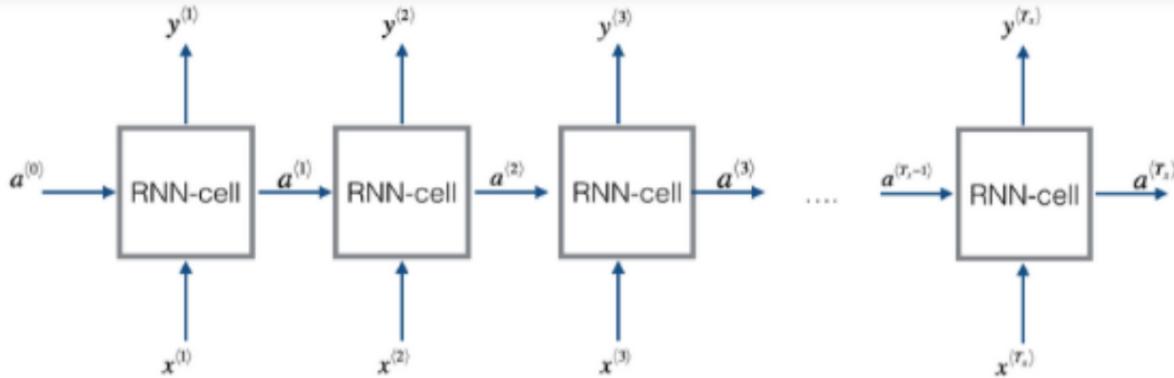


Figure 1: Basic RNN model

Dimensions of input x

Input with n_x number of units

- For a single time step of a single input example, $x^{(i)(t)}$ is a one-dimensional input vector
- Using language as an example, a language with a 5000-word vocabulary could be one-hot encoded into a vector that has 5000 units. So $x^{(i)(t)}$ would have the shape (5000,)
- The notation n_x is used here to denote the number of units in a single time step of a single training example

Time steps of size T_x

- A recurrent neural network has multiple time steps, which you'll index with t .
- In the lessons, you saw a single training example $x^{(i)}$ consisting of multiple time steps T_x . In this notebook, T_x will denote the number of timesteps in the longest sequence.

Batches of size m

- Let's say we have mini-batches, each with 20 training examples
- To benefit from vectorization, you'll stack 20 columns of $x^{(i)}$ examples
- For example, this tensor has the shape (5000,20,10)
- You'll use m to denote the number of training examples
- So, the shape of a mini-batch is (n_x, m, T_x)

3D Tensor of shape (n_x, m, T_x)

- The 3-dimensional tensor x of shape (n_x, m, T_x) represents the input x that is fed into the RNN

Taking a 2D slice for each time step: $x^{(t)}$

- At each time step, you'll use a mini-batch of training examples (not just a single example)
- So, for each time step t , you'll use a 2D slice of shape (n_x, m)
- This 2D slice is referred to as $x^{(t)}$. The variable name in the code is `xt`.

Definition of hidden state a

- The activation $a^{(t)}$ that is passed to the RNN from one time step to another is called a "hidden state."

Dimensions of hidden state a

- Similar to the input tensor x , the hidden state for a single training example is a vector of length n_a
- If you include a mini-batch of m training examples, the shape of a mini-batch is (n_a, m)
- When you include the time step dimension, the shape of the hidden state is (n_a, m, T_x)
- You'll loop through the time steps with index t , and work with a 2D slice of the 3D tensor
- This 2D slice is referred to as $a^{(t)}$
- In the code, the variable names used are either `a_prev` or `a_next`, depending on the function being implemented
- The shape of this 2D slice is (n_a, m)

Dimensions of prediction \hat{y}

- Similar to the inputs and hidden states, \hat{y} is a 3D tensor of shape (n_y, m, T_y)
 - n_y : number of units in the vector representing the prediction
 - m : number of examples in a mini-batch
 - T_y : number of time steps in the prediction
- For a single time step t , a 2D slice $\hat{y}^{(t)}$ has shape (n_y, m)
- In the code, the variable names are:
 - `y_pred` : \hat{y}
 - `yt_pred` : $\hat{y}^{(t)}$

Here's how you can implement an RNN:

Steps:

- Implement the calculations needed for one time step of the RNN.
- Implement a loop over T_x time steps in order to process all the inputs, one at a time.

1.1 - RNN Cell

You can think of the recurrent neural network as the repeated use of a single cell. First, you'll implement the computations for a single time step. The following figure describes the operations for a single time step of an RNN cell:

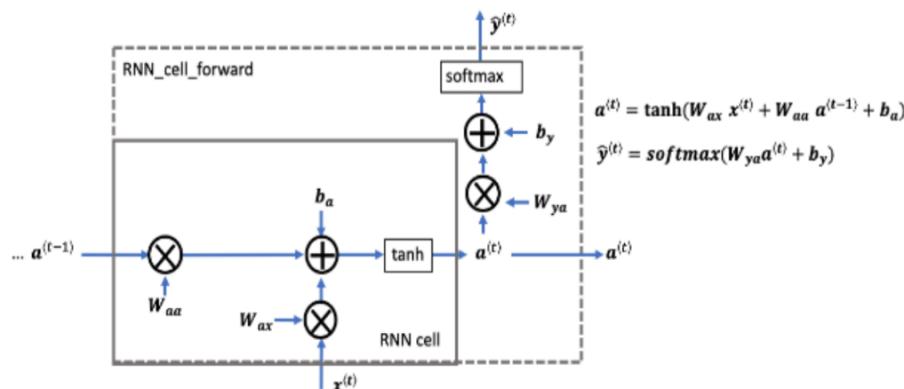


Figure 2: Basic RNN cell. Takes as input $x^{(t)}$ (current input) and $a^{(t-1)}$ (previous hidden state containing information from the past), and outputs $a^{(t)}$ which is given to the next RNN cell and also used to predict $\hat{y}^{(t)}$

1.2 - RNN Forward Pass

- A recurrent neural network (RNN) is a repetition of the RNN cell that you've just built.
 - If your input sequence of data is 10 time steps long, then you will re-use the RNN cell 10 times
- Each cell takes two inputs at each time step:
 - $a^{(t-1)}$: The hidden state from the previous cell
 - $x^{(t)}$: The current time step's input data
- It has two outputs at each time step:
 - A hidden state ($a^{(t)}$)
 - A prediction ($\hat{y}^{(t)}$)
- The weights and biases (W_{aa}, b_a, W_{ax}, b_y) are re-used each time step
 - They are maintained between calls to `rnn_cell_forward` in the 'parameters' dictionary

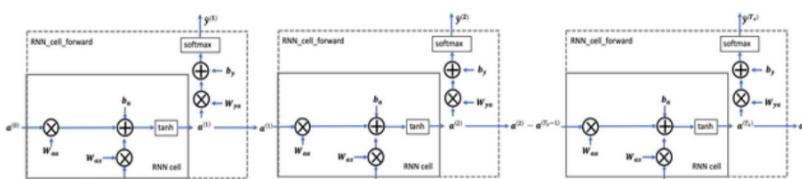


Figure 3: Basic RNN. The input sequence $x = (x^{(1)}, x^{(2)}, \dots, x^{(T_x)})$ is carried over T_x time steps. The network outputs $y = (y^{(1)}, y^{(2)}, \dots, y^{(T_x)})$.

Exercise 2 - rnn_forward

Implement the forward propagation of the RNN described in Figure 3.

Instructions:

- Create a 3D array of zeros, a of shape (n_a, m, T_x) that will store all the hidden states computed by the RNN
- Create a 3D array of zeros, \hat{y} , of shape (n_y, m, T_x) that will store the predictions
 - Note that in this case, $T_y = T_x$ (the prediction and input have the same number of time steps)
- Initialize the 2D hidden state `a_next` by setting it equal to the initial hidden state, a_0
- At each time step t :
 - Get $x^{(t)}$, which is a 2D slice of x for a single time step t
 - $x^{(t)}$ has shape (n_x, m)
 - x has shape (n_x, m, T_x)
 - Update the 2D hidden state $a^{(t)}$ (variable name `a_next`), the prediction $\hat{y}^{(t)}$ and the cache by running `rnn_cell_forward`
 - $a^{(t)}$ has shape (n_a, m)
 - Store the 2D hidden state in the 3D tensor a , at the t^{th} position
 - a has shape (n_a, m, T_x)
 - Store the 2D $\hat{y}^{(t)}$ prediction (variable name `yt_pred`) in the 3D tensor \hat{y}_{pred} at the t^{th} position
 - $\hat{y}^{(t)}$ has shape (n_y, m)
 - \hat{y} has shape (n_y, m, T_x)
 - Append the cache to the list of caches
- Return the 3D tensor a and \hat{y} , as well as the list of caches

