



## Assignment 4 Building Recurrent Neural Network

### 1 Objective

- Building forward and backward directions of a basic RNN cell.

### 2 Problem Statement

Recurrent Neural Networks (RNN) are very effective for Natural Language Processing and other sequence tasks because they have "memory". They can read inputs  $x^{(t)}$  (such as words) one at a time, and remember some information/context through the hidden layer activations that get passed from one time-step to the next. This allows a uni-directional RNN to take information from the past to process later inputs. A bidirection RNN can take context from both the past and the future.

These are the notations that we will use in the assignment:

- Superscript  $[l]$  denotes an object associated with the  $l^{th}$  layer.  
Example:  $a^{[4]}$  is the 4<sup>th</sup> layer activation.  $W^{[5]}$  and  $b^{[5]}$  are the 5<sup>th</sup> layer parameters.
- Superscript  $(i)$  denotes an object associated with the  $i^{th}$  example.  
Example:  $x^{(i)}$  is the  $i^{th}$  training example input.
- Superscript  $\langle t \rangle$  denotes an object at the  $t^{th}$  time-step.  
Example:  $x^{\langle t \rangle}$  is the input  $x$  at the  $t^{th}$  time-step.  $x^{(i)\langle t \rangle}$  is the input at the  $t^{th}$  timestep of example  $i$ .
- Lowerscript  $i$  denotes the  $i^{th}$  entry of a vector. Example:  $a_i^{[l]}$  denotes the  $i^{th}$  entry of the activations in layer  $l$ .

## 2.1 Forward propagation for the basic Recurrent Neural Network

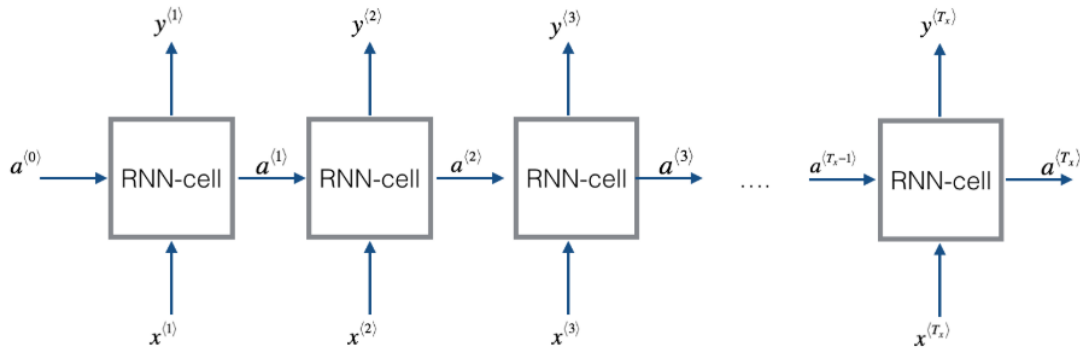


Figure 1

Here's how you can implement an RNN:

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

### 2.1.1 RNN cell

A Recurrent neural network can be seen as the repetition of a single cell. You are first going to implement the computations for a single time-step. The following figure describes the operations for a single time-step of an RNN cell.

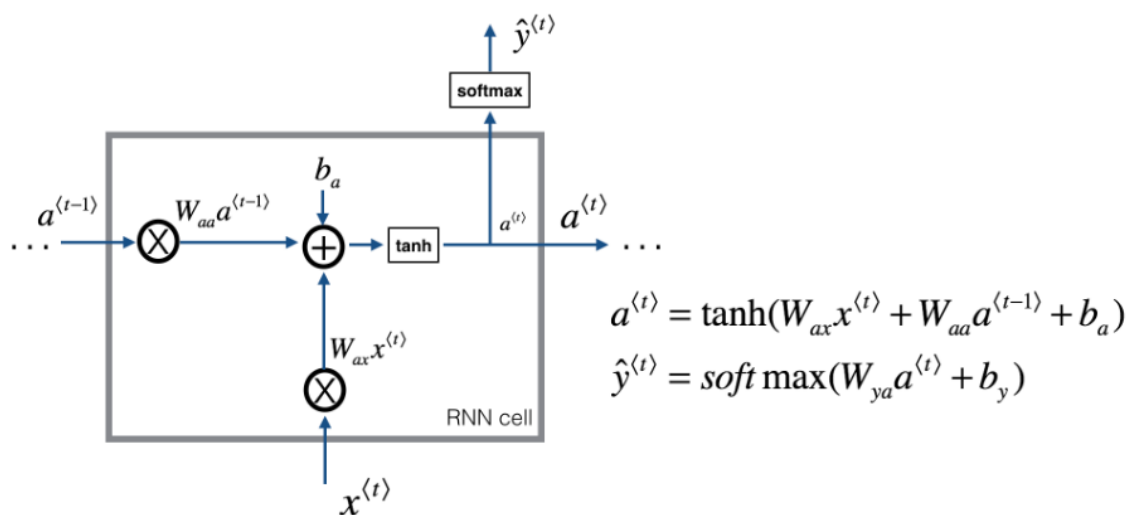


Figure 2

### 2.1.2 RNN forward pass

You can see an RNN as the repetition of the cell you've just built. If your input sequence of data is carried over 10 time steps, then you will copy the RNN cell 10 times. Each cell takes as input the hidden state from the previous cell ( $a^{(t-1)}$ ) and the current time-step's input data ( $x^{(t)}$ ). It outputs a hidden state ( $a^{(t)}$ ) and a prediction ( $y^{(t)}$ ) for this time-step.

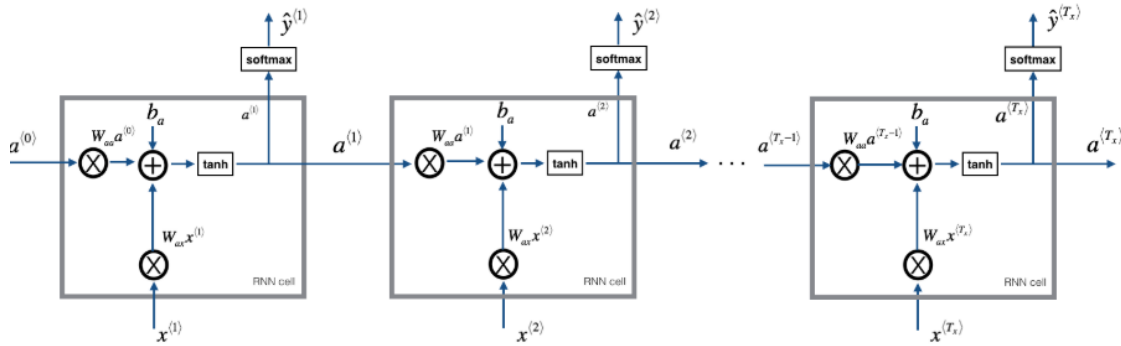


Figure 3

## 3 Long Short-Term Memory (LSTM) network

This following figure shows the operations of an LSTM-cell

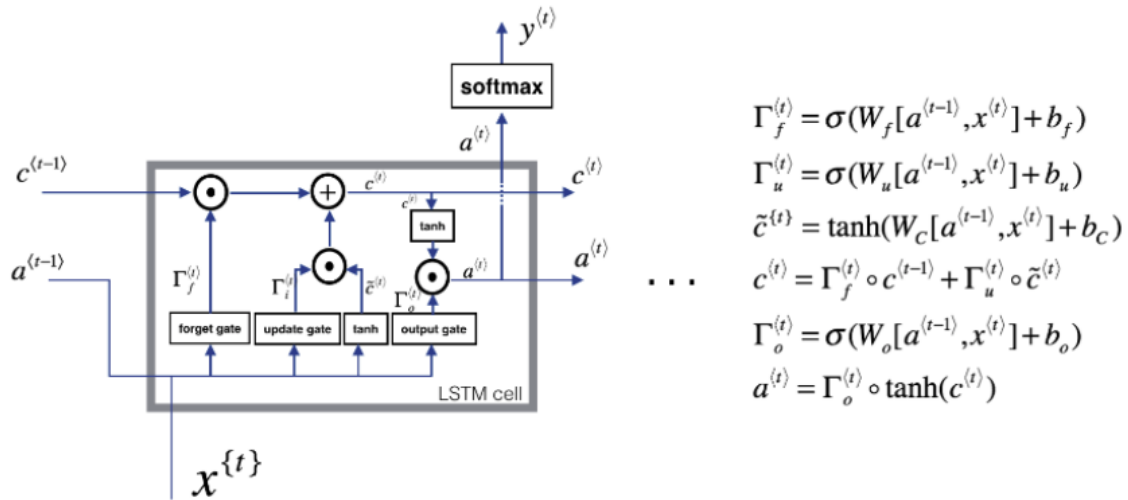


Figure 4

Similar to the RNN example above, you will start by implementing the LSTM cell for a single time-step. Then you can iteratively call it from inside a for-loop to have it process an



input with  $T_x$  time-steps.

## About the gates

### -Forget gate

For the sake of this illustration, let's assume we are reading words in a piece of text, and want to use an LSTM to keep track of grammatical structures, such as whether the subject is singular or plural. If the subject changes from a singular word to a plural word, we need to find a way to get rid of our previously stored memory value of the singular/plural state. In an LSTM, the forget gate lets us do this:

$$\Gamma_f^{(t)} = \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f)$$

Here,  $W_f$  are weights that govern the forget gate's behavior. We concatenate  $[a^{(t-1)}, x^{(t)}]$  and multiply by  $W_f$ . The equation above results in a vector  $\Gamma_f^{(t)}$  with values between 0 and 1. This forget gate vector will be multiplied element-wise by the previous cell state  $c^{(t-1)}$ . So if one of the values of  $\Gamma_f^{(t)}$  is 0 (or close to 0) then it means that the LSTM should remove that piece of information (e.g. the singular subject) in the corresponding component of  $c^{(t-1)}$ . If one of the values is 1, then it will keep the information.

### - Update gate

Once we forget that the subject being discussed is singular, we need to find a way to update it to reflect that the new subject is now plural. Here is the formula for the update gate:

$$\Gamma_u^{(t)} = \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u)$$

Similar to the forget gate, here  $\Gamma_u^{(t)}$  is again a vector of values between 0 and 1. This will be multiplied element-wise with  $\tilde{c}^{(t)}$ , in order to compute  $c^{(t)}$ .

**- Updating the cell** To update the new subject we need to create a new vector of numbers that we can add to our previous cell state. The equation we use is:

$$\tilde{c}^{(t)} = \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c)$$

Finally, the new cell state is:

0

$$c^{(t)} = \Gamma_f^{(t)} * c^{(t-1)} + \Gamma_u^{(t)} * \tilde{c}^{(t)}$$

**- Output gate** To decide which outputs we will use, we will use the following two formulas:

$$\Gamma_o^{(t)} = \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o)$$

$$a^{(t)} = \Gamma_o^{(t)} * \tanh(c^{(t)})$$

Where in equation 5 you decide what to output using a sigmoid function and in equation 6 you multiply that by the tanh of the previous state.

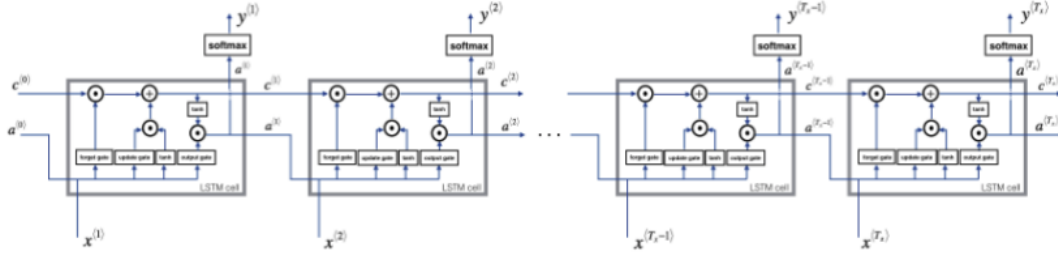


Figure 5

## 4 Backpropagation in recurrent neural networks

### 4.1 Basic RNN backward pass

Here we will compute the backward pass for the basic RNN-cell.

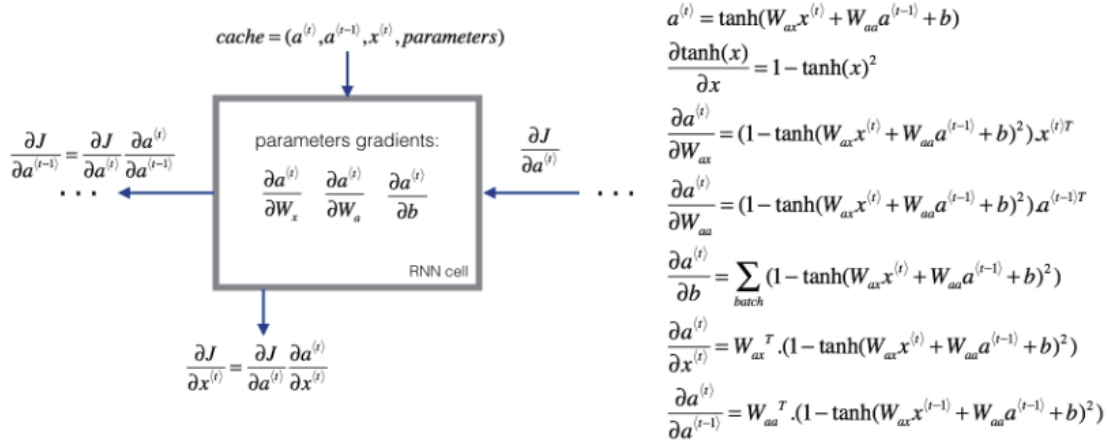


Figure 6

Deriving the one step backward functions: To compute the rnn\_cell.backward you need to compute the following equations.

The derivative of tanh is  $1 - \tanh(x)^2$ . Note that:  $\sec(x)^2 = 1 - \tanh(x)^2$

Similarly for  $\frac{\partial a^{(t)}}{\partial W_{ax}}$ ,  $\frac{\partial a^{(t)}}{\partial W_{aa}}$ ,  $\frac{\partial a^{(t)}}{\partial b}$ , the derivative of tanh(u) is  $(1 - \tanh(u)^2)du$ .



The final two equations also follow the same rule and are derived using the tanh derivative. Note that the arrangement is done in a way to get the same dimensions to match.

## 5 Requirements

- Implement the RNN-cell described in Figure 2.
  - Compute the hidden state with tanh activation:  $a^{(t)} = \tanh(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a)$ .
  - Using your new hidden state  $a^{(t)}$ , compute the prediction  $\hat{y}^{(t)} = \text{softmax}(W_{ya}a^{(t)} + b_y)$ . We provided you a function: softmax.
  - Store  $(a^{(t)}, a^{(t-1)}, x^{(t)}, \text{parameters})$  in cache
  - Return  $a^{(t)}$ ,  $y^{(t)}$  and cache

We will vectorize over  $m$  examples. Thus,  $x^{(t)}$  will have dimension  $(n_x, m)$ , and  $a^{(t)}$  will have dimension  $(n_a, m)$ .

Expected Output:

|                          |   |
|--------------------------|---|
| <b>**a_next[4]**:</b>    | [ 0.59584544 0.18141802 0.61311866 0.99808218 0.85016201 0.99980978 -0.18887155 0.99815551<br>0.6531151 0.82872037] |
| <b>**a_next.shape**:</b> | (5, 10)   |
| <b>**yt[1]**:</b>        | [ 0.9888161 0.01682021 0.21140899 0.36817467 0.98988387 0.88945212 0.36920224 0.9966312 0.9982559<br>0.17746526]    |
| <b>**yt.shape**:</b>     | (2, 10)   |

- Code the forward propagation of the RNN described in Figure 3.
  - Create a vector of zeros ( $a$ ) that will store all the hidden states computed by the RNN.
  - Initialize the “next” hidden state as  $a_0$  (initial hidden state).
  - Start looping over each time step, your incremental index is  $t$  :
    - \* Update the “next” hidden state and the cache by running rnn\_cell\_forward
    - \* Store the “next” hidden state in  $a$  ( $t^{\text{th}}$  position)
    - \* Store the prediction in  $y$
    - \* Add the cache to the list of caches
  - Return  $a$ ,  $y$  and caches.



Expected Output:

|                            |   |
|----------------------------|---|
| <b>**a[4][1]**:</b>        | <b>[-0.99999375 0.77911235 -0.99861469 -0.99833267]</b> |
| <b>**a.shape**:</b>        | <b>(5, 10, 4)</b>                                       |
| <b>**y[1][3]**:</b>        | <b>[ 0.79560373 0.86224861 0.11118257 0.81515947]</b>   |
| <b>**y.shape**:</b>        | <b>(2, 10, 4)</b>                                       |
| <b>**cache[1][1][3]**:</b> | <b>[-1.1425182 -0.34934272 -0.20889423 0.58662319]</b>  |
| <b>**len(cache)**:</b>     | <b>2</b>  |

- Implement the LSTM cell described in the Figure 4.

- Concatenate  $a^{(t-1)}$  and  $x^{(t)}$  in a single matrix:  $concat = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$
  - Compute all the formulas 1-6. You can use `sigmoid()` (provided) and `np.tanh()`.
  - Compute the prediction  $y^{(t)}$ . You can use `softmax()` (provided).
- Expected Output:

|                          |   |
|--------------------------|---|
| <b>**a_next[4]**:</b>    | <b>[-0.66408471 0.0036921 0.02088357 0.22834167 -0.85575339 0.00138482 0.76566531 0.34631421 -0.00215674 0.43827275]</b>    |
| <b>**a_next.shape**:</b> | <b>(5, 10)</b>  |
| <b>**c_next[2]**:</b>    | <b>[ 0.63267805 1.00570849 0.35504474 0.20690913 -1.64566718 0.11832942 0.76449811 -0.0981561 -0.74348425 -0.26810932]</b>  |
| <b>**c_next.shape**:</b> | <b>(5, 10)</b>  |
| <b>**yt[1]**:</b>        | <b>[ 0.79913913 0.15986619 0.22412122 0.15606108 0.97057211 0.31146381 0.00943007 0.12666353 0.39380172 0.07828381]</b>     |
| <b>**yt.shape**:</b>     | <b>(2, 10)</b>  |
| <b>**cache[1][3]**:</b>  | <b>[-0.16263996 1.03729328 0.72938082 -0.54101719 0.02752074 -0.30821874 0.07651101 -1.03752894 1.41219977 -0.37647422]</b> |
| <b>**len(cache)**:</b>   | <b>10</b>   |



- Implement `lstm_forward()` to run an LSTM over  $T_x$  time-steps, Note that  $c^{(0)}$  is initialized with zeros.

Expected Output:

|                                    |   |
|------------------------------------|---|
| <code>**a[4][3][6]** =</code>      | 0.172117767533  |
| <code>**a.shape** =</code>         | (5, 10, 7)  |
| <code>**y[1][4][3]** =</code>      | 0.95087346185   |
| <code>**y.shape** =</code>         | (2, 10, 7)  |
| <code>**caches[1][1][1]** =</code> | [ 0.82797464 0.23009474 0.76201118 -0.22232814 -0.20075807 0.18656139 0.41005165] |
| <code>**c[1][2][1]** =</code>      | -0.855544916718   |
| <code>**len(caches)** =</code>     | 2   |

- Implement `rnn_cell_backward()` to run an LSTM over  $T_x$  time-steps, Note that  $c^{(0)}$  is initialized with zeros.

Expected Output:

|   |                 |
|---|-----------------|
| <code>**gradients["dxt"][1][2]** =</code>     | -0.460564103059 |
| <code>**gradients["dxt"].shape** =</code>     | (3, 10)         |
| <code>**gradients["da_prev"][2][3]** =</code> | 0.0842968653807 |
| <code>**gradients["da_prev"].shape** =</code> | (5, 10)         |
| <code>**gradients["dWax"][3][1]** =</code>    | 0.393081873922  |
| <code>**gradients["dWax"].shape** =</code>    | (5, 3)          |
| <code>**gradients["dWaa"][1][2]** =</code>    | -0.28483955787  |
| <code>**gradients["dWaa"].shape** =</code>    | (5, 5)          |
| <code>**gradients["dba"][4]** =</code>        | [ 0.80517166]   |
| <code>**gradients["dba"].shape** =</code>     | (5, 1)          |

- Implement the `rnn_backward` function. Initialize the return variables with zeros first and then loop through all the time steps while calling the `rnn_cell_backward` at each time timestep, update the other variables accordingly.





Expected Output:

|  |  |
|--|--|
| <code>**gradients["dx"][1][2]** =</code>   | <code>[-2.07101689 -0.59255627 0.02466855 0.01483317]</code> |
| <code>**gradients["dx"].shape** =</code>   | <code>(3, 10, 4)</code>                                      |
| <code>**gradients["da0"][2][3]** =</code>  | <code>-0.314942375127</code>                                 |
| <code>**gradients["da0"].shape** =</code>  | <code>(5, 10)</code>   |
| <code>**gradients["dWax"][3][1]** =</code> | <code>11.2641044965</code>                                   |
| <code>**gradients["dWax"].shape** =</code> | <code>(5, 3)</code>  |
| <code>**gradients["dWaa"][1][2]** =</code> | <code>2.30333312658</code>                                   |
| <code>**gradients["dWaa"].shape** =</code> | <code>(5, 5)</code>  |
| <code>**gradients["dba"][4]** =</code>     | <code>[-0.74747722]</code>                                   |
| <code>**gradients["dba"].shape** =</code>  | <code>(5, 1)</code>  |

## 6 BONUS: Sentiment Analysis

Sentiment analysis is one of the most common applications in Natural Language processing. So here we are, we will train a classifier for movie reviews in IMDB data set, using TF.Keras Recurrent Neural Networks.

1. Load the IMDB data from Keras
2. Build the model.
3. Decide the percentage of the training data that you will use as validation
4. Train the model, then evaluate it.
5. Report the best accuracy you got in each of the 3 considered alternatives.

## 7 Notes

- Parts of this assignment are based on Andrew Ng Stanford deep learning course.
- The dataset files and the starter code for the problems can be found in the resources section in Piazza.



- Cheating will be severely penalized (for both parties). So, it is better to deliver nothing than deliver a copy. Any online resources used must be clearly identified.