

## LSTM HW6

1. *Proof.* From the definition given by PyTorch, we know that the simple RNN with sigmoid nonlinearity updates its hidden state  $h_t$  using the following equation:

$$h_t = \sigma(x_t W_{ih}^T + h_{t-1} W_{hh}^T + b_{ih} + b_{hh})$$

where  $\sigma$  is the sigmoid function. Also from PyTorch, we know that the LSTM is defined as follow:

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ h_t &= o_t \odot \tanh(c_t). \end{aligned}$$

Notice that  $o_t$  is the gate with sigmoid non-linearity. Hence, we can show that the simple RNN with sigmoid non-linearity is a subset of LSTM letting  $\tanh(c_t) \approx 1$ . We know that  $\lim_{c_t \rightarrow \infty} \tanh(c_t) = 1$ . This means that we need  $c_t$  to be very large.

$f_t$  Notice that since  $f_t$  uses the sigmoid activation function,  $0 < f_t < 1$ .  $f_t$  controls how much of the previous cell state  $c_{t-1}$  is retained. Then to ensure large increases in  $c_t$ , we need to set  $f_t \rightarrow 1$ .

$c_0$  To ensure  $c_t$  diverges to infinity, we can set the initial state with large positive values.

$i_t$  Notice that  $i_t$  also involves the sigmoid activation function, which means that  $0 < i_t < 1$ .  $i_t$  controls how much of the new information is written to the cell state. Then we can set  $i_t \rightarrow 1$  to ensure  $c_t$  goes to infinity.

$g_t$  Since  $g_t$  involve  $\tanh$ ,  $-1 < g_t < 1$ . Then we can set  $g_t \rightarrow 1$  for  $c_t$  to approach to infinity.

Now we can rewrite the recurrence function as

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ &= f_t^t c_0 + t \cdot 1 \cdot 1 \\ &= f_t^t c_0 + t. \end{aligned}$$

Notice that  $\lim_{t \rightarrow \infty} c_t = \infty$ . Also, we mentioned that  $c_0$  can be initialized with large values, then  $\tanh(c_t) \rightarrow 1$ . Then,

$$h_t \rightarrow o_t \odot 1 = \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}),$$

which is the recurrence equation of a simple RNN with sigmoid nonlinearity. Hence, the simple RNN with sigmoid nonlinearity is a subset of the LSTM.

For the second part of the question, we can set all weights  $W = 0$  with a large bias. Then,

$$\begin{aligned} f_t &= i_t = \sigma(0 + \text{large bias}) \rightarrow 1 \\ g_t &= \tanh(0 + \text{large bias}) \rightarrow 1. \end{aligned}$$

Then these conditions will ensure the LSTM to mimic the RNN by what we have shown above.  $\square$

## 2. One-hot Encoding

According to the alphabetical order, we set

$$\begin{aligned}\text{"bad"} &= [1, 0, 0, 0] \\ \text{"good"} &= [0, 1, 0, 0] \\ \text{"not"} &= [0, 0, 1, 0] \\ \text{"uh"} &= [0, 0, 0, 1].\end{aligned}$$

### General Workflow

Notice the recurrence relationship of  $c_t$  is

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t.$$

We want to delay the sentiment score by 1. That is,  $c_t$  contains the cumulated sentiment score at  $t - 1$ ,  $c_{t-1}$  contains the cumulated sentiment score at  $t - 2$ , and  $g_t$  calculates the sentiment score at time  $t - 1$  by utilizing the hidden state  $h_{t-1}$ . Let's then consider each gate separately.

### Gate Specification

**Input  $i_t$**  We want  $i_t = \mathbf{1}_{1 \times 6}$  for the following reason:

The equation we are given is

$$\sum \text{"good"} - \sum \text{"bad"} - 2 \sum \text{"not good"} + 2 \sum \text{"not bad"}.$$

This means that in each iteration, we are either adding 1, subtracting 1 or doing nothing (when the current word is "not" or "uh"). Thus, we need to completely keep the updated information, namely,  $i_t = [1, 1, 1, 1, 1, 1]^T$ . To do this, we just need to set

$$\begin{aligned}W_{ii} &= \mathbf{0}_{6 \times 4} \\ W_{hi} &= \mathbf{0}_{6 \times 6} \\ b_{ii} &= \mathbf{0}_{1 \times 6} \\ b_{hi} &= \mathbf{100}_{1 \times 6}.\end{aligned}$$

Then, we have

$$\begin{aligned}i_t &= \sigma(0 + b_{ii} + 0 + b_{hi}) \\ &= \sigma([100, 100, 100, 100, 100, 100]^T) \\ &\approx [1, 1, 1, 1, 1, 1]^T.\end{aligned}$$

**Forget  $f_t$**  We want  $f_t = [0, 0, 0, 0, 0, 1]^T$  for the following reason:

In the forget gate, since we only care about the sentiment score, we only need to set the element corresponding to the sentiment score to 1. In our case, we set the last element to be the sentiment score. To achieve this, we just need to set

$$\begin{aligned}W_{if} &= \mathbf{0}_{6 \times 4} \\ W_{hf} &= \mathbf{0}_{6 \times 6} \\ b_{if} &= [-100, -100, -100, -100, -100, 100]^T \\ b_{hf} &= \mathbf{0}_{1 \times 6}.\end{aligned}$$

Then, we have

$$\begin{aligned}f_t &= \sigma(0 + b_{if} + 0 + 0) \\ &= \sigma([-100, -100, -100, -100, -100, 100]^T) \\ &\approx [0, 0, 0, 0, 0, 1].\end{aligned}$$

**Cell  $g_t$**  We want  $g_t = [1, 1, 1, 1, 1, \text{delayed increment score}]^T$ .

The reason that we want the first five to be 1 is as follows: The update equations are:

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ h_t &= o_t \odot \tanh(c_t). \end{aligned}$$

We want to achieve  $h_t = [I_{\text{bad}}, I_{\text{good}}, I_{\text{not}}, I_{\text{not bad}}, I_{\text{not good}}, 0]^T$ , which stores the flags. Notice that the first five elements of  $f_t \odot c_{t-1}$  are zeros. If we set the first five elements of  $g_t$  to zero as well, then the first five elements of  $c_t$  will remain zero, since both terms in the update equation contribute zeros for these elements. Consequently,  $\tanh(c_t) = 0$  for the first five elements, causing  $h_t$  to lack any meaningful information in these positions. To store the flag information in  $h_t$ , we therefore need to set the first five elements of  $g_t$  to ones. This ensures that the input gate  $i_t \odot g_t$  can directly influence  $c_t$  for these elements, allowing the desired flag indicators to be reflected in  $h_t$ .

Now, let's propagate the weights and biases. Notice that

$$g_t = \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}),$$

and we only want to use this gate to determine the sentiment score at  $t - 1$  through the hidden state  $h_{t-1}$ . To achieve this, since we do not need to handle  $x_t$  at this point, we can set  $W_{ig} = \mathbf{0}_{6 \times 4}$  and  $b_{ig} = \mathbf{0}_{1 \times 6}$ . To determine the sentiment score, according to the equation, we have

$$\begin{aligned} \text{bad} &= -1 \\ \text{good} &= +1 \\ \text{not bad} &= 2 \cdot 1 - 1 = +1 \text{ (since "bad" alone will also be counted)} \\ \text{not good} &= -2 \cdot 1 + 1 = -1 \text{ (since "good" alone will also be counted)}. \end{aligned}$$

Thus, we can design  $W_{hg}$  as follows:

$$W_{hg} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -100 & 100 & 0 & 100 & -100 & 0 \end{bmatrix}.$$

Recall that we design our hidden state matrix to be flags of the words, that is

$$h_{t-1} = [I_{\text{bad}}, I_{\text{good}}, I_{\text{not}}, I_{\text{not bad}}, I_{\text{not good}}, 0]^T.$$

Then,

$$(W_{hg}h_{t-1})_6 = -100I_{\text{bad}} + 100I_{\text{good}} + 100I_{\text{not bad}} - 100I_{\text{not good}},$$

which serves the function of determining sentiment scores. The first five elements of  $W_{hg}h_{t-1}$  are zeros now. We then set

$$b_{ig} = [100, 100, 100, 100, 100, 0]^T.$$

Then,

$$\begin{aligned} g_t &= \tanh(0 + 0 + W_{hg}h_{t-1} + b_{hg}) \\ &= \tanh([100, 100, 100, 100, 100, (W_{hg}h_{t-1})_6]^T) \\ &= [1, 1, 1, 1, 1, \text{delayed increment score}]^T \end{aligned}$$

where  $\text{delayed increment score} = \tanh((W_{hg}h_{t-1})_6) \in \{-1, 0, 1\}$ .

**Cell State  $c_t$**  Using the settings from above, it follows that our cell state will be

$$c_t = [1, 1, 1, 1, 1, \text{delayed cumulated score}]^T.$$

**Output  $o_t$**  We want  $o_t$  to store flags of the current word given the delayed information stored in  $h_{t-1}$ . That is, we want

$$o_t = [I_{\text{bad}}, I_{\text{good}}, I_{\text{not}}, I_{\text{not bad}}, I_{\text{not good}}, 0]^T.$$

To do so, we need  $W_{io}x_t$  to handle the current word, and  $W_{ho}h_{t-1}$  to determine if the previous word is "not". Then we set

$$W_{io} = \begin{bmatrix} 100 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 \\ 0 & 0 & 100 & 0 \\ 100 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

and

$$W_{ho} = \begin{bmatrix} 0 & 0 & -100 & 0 & 0 & 0 \\ 0 & 0 & -100 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 100 & 0 & 0 & 0 \\ 0 & 0 & 100 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Then,  $W_{io}x_t = [100I_{\text{bad}}, 100I_{\text{good}}, 100I_{\text{not}}, 100I_{\text{bad}}, 100I_{\text{good}}, 0]^T$ , and  $W_{ho}h_{t-1} = [-100I'_{\text{not}}, -100I'_{\text{not}}, 0, 100I'_{\text{not}}, 100I'_{\text{not}}, 0]^T$  where  $I'_{\text{not}}$  is the "not" flag at  $t-1$ . Then, we set

$$b_{ho} = [-50, -50, -50, -150, -150, -50]^T.$$

The bias values are chosen this way to selectively activate elements of  $o_t$  based on specific inputs. For example, consider the phrase "not bad," where we want  $o_t = [0, 0, 0, 1, 0, 0]^T$ . Given that:

$$W_{io}x_t = [100, 0, 0, 100, 0, 0]^T,$$

$$W_{ho}h_{t-1} = [-100, -100, 0, 100, 100, 0]^T,$$

we have

$$W_{io}x_t + W_{ho}h_{t-1} = [0, -100, 0, 200, 100, 0]^T.$$

Since  $\sigma(0) = 0.5$ , without additional scaling, multiple elements might be activated or stay in the intermediate state. To ensure only the intended element activates, we add a large negative bias to reduce the unwanted elements' values, thereby ensuring only the correct position passes through the sigmoid function. These bias values were obtained by iterative testing across various cases to achieve selective activation in  $o_t$ .

**Hidden  $h_t$**  Lastly, we update the hidden state at last using the function  $h_t = o_t \odot \tanh(c_t)$ . As the last state getting updated,  $h_t$  stores the information of the previous step.

## Result Validation

We use Python to validate if our configuration is correct

```

1  ### Author: Yuxuan Zhang
2  ### Date: 10/24/2024
3  import numpy as np
4  def sigmoid(x):
5      return 1 / (1 + np.exp(-x))
6
7  def tanh(x):
8      return np.tanh(x)
9
10 def lstm(sequence):
11     # Define the one-hot encoding
12     word_dict = {
13         "bad": np.array([1, 0, 0, 0]),
14         "good": np.array([0, 1, 0, 0]),
15         "not": np.array([0, 0, 1, 0]),
16         "uh": np.array([0, 0, 0, 1])
17     }
18
19     # Initialize h_0 and c_0
20     h_t = np.zeros(6)
21     c_t = np.zeros(6)
22
23     # For simplicity, set i_t and f_t to constants
24     i_t = np.ones(6)
25     f_t = np.array([0, 0, 0, 0, 0, 1])
26
27     # Cell Gate matrices
28     W_ig = np.zeros((6, 4))
29     W_hg = np.array([
30         [0, 0, 0, 0, 0, 0],
31         [0, 0, 0, 0, 0, 0],
32         [0, 0, 0, 0, 0, 0],
33         [0, 0, 0, 0, 0, 0],
34         [0, 0, 0, 0, 0, 0],
35         [-100, 100, 0, 100, -100, 0]
36     ])
37
38     b_ig = np.array([100, 100, 100, 100, 100, 0])
39     b_hg = np.zeros(6)
40
41     # Output gate matrices
42     W_io = np.array([
43         [100, 0, 0, 0],
44         [0, 100, 0, 0],
45         [0, 0, 100, 0],
46         [100, 0, 0, 0],
47         [0, 100, 0, 0],
48         [0, 0, 0, 0]
49     ])
50     W_ho = np.array([
51         [0, 0, -100, 0, 0, 0],
52         [0, 0, -100, 0, 0, 0],
53         [0, 0, 0, 0, 0, 0],
54         [0, 0, 100, 0, 0, 0],
55         [0, 0, 100, 0, 0, 0],
56         [0, 0, 0, 0, 0, 0]
57     ])
58
59     b_io = np.zeros(6)
60     b_ho = np.array([-50, -50, -50, -150, -150, -50] )

```

```

61
62     for word in sequence:
63         print(f"Processing word: {word}")
64
65         x_t = word_dict[word]
66         g_t = tanh(np.dot(W_ig, x_t) + b_ig + np.dot(W_hg, h_t) + b_hg)
67         c_t = f_t * c_t + i_t * g_t
68         o_t = sigmoid(np.dot(W_io, x_t) + b_io + np.dot(W_ho, h_t) + b_ho)
69         h_t = o_t * tanh(c_t)
70
71         print("Updated delayed total sentiment score:", c_t[-1])
72         print("-" * 50)
73
74 def main():
75     sequence = ["uh", "good", "good", "not", "not", "bad", "bad", "uh"]
76     lstm(sequence)
77
78 if __name__ == "__main__":
79     main()

```

d:/NLP-Course-Project/LSTM/LSTM.py

Processing word: uh  
Updated delayed total sentiment score: 0.0

-----  
Processing word: good  
Updated delayed total sentiment score: 0.0

-----  
Processing word: good  
Updated delayed total sentiment score: 1.0

-----  
Processing word: not  
Updated delayed total sentiment score: 2.0

-----  
Processing word: not  
Updated delayed total sentiment score: 2.0

-----  
Processing word: bad  
Updated delayed total sentiment score: 2.0

-----  
Processing word: bad  
Updated delayed total sentiment score: 3.0

-----  
Processing word: uh  
Updated delayed total sentiment score: 2.0

-----  
The sentence we used is "uh good good not not bad bad uh." Note that we added "uh" to the end of the sequence as a placeholder for the delayed sentiment score calculation. In this sequence, there are 2 occurrences of "good," 2 occurrences of "bad," and 1 occurrence of "not bad." Using the sentiment equation, we find the total sentiment score to be

$$1 \cdot 2 - 1 \cdot 2 + 1 \cdot 2 = 2.$$

This result matches the output from our code, which validates our configuration choices.