

LSTM

Table of Contents

1. LSTM	1
1.1. Mathematical Formulation	1
2. Code Implementation	2

1. LSTM

1.1. Mathematical Formulation

When we train RNN, we usually face the problem of vanishing gradient or exploding gradient.

To solve this problem, we can use LSTM, whose architecture use simple addition and multiplication to update the memory.

LSTM introduces some new components compared to RNN:

Definition 1.1

- Memory Cell C_t : The memory cell is a vector that stores the long-term memory. It updates at every time step.
- Gate:
 - Input Gate i_t : Control what information should be added to the memory.
 - Forget Gate f_t : Control what information should be forgotten from the memory.
 - Output Gate o_t : Control what information should be output to the next memory cell C_{t+1} .
- Hidden State h_t : Just like RNN, LSTM has h_t , but the hidden state is a vector that stores the short-term memory.

Before introducing the forward process, we need to introduce the hadamard product.

Definition 1.2

The hadamard product of two vectors $a \in \mathbb{R}^n$ and $b \in \mathbb{R}^n$ is defined as:

$$a \odot b = [a_1 b_1, a_2 b_2, \dots, a_n b_n]$$

For a LSTM of $[0, T]$ time steps, the whole forward process can be described as:

Definition 1.3

At timestep t :

1. Gate update:

- Forget Gate: $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$
- Input Gate: $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$, $\tilde{C}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_C)$
- Output Gate: $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$

-
2. Memory Cell update: $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$
 3. Hidden State update: $h_t = o_t \cdot \tanh(C_t)$
 4. (Optional) Output: $y_t = f(W_y h_t + b_y)$

Let's do a **dimension analysis**:

Example 1.4

For an input sequence x of size $N_{\text{input}} = D$, the dimension of different parts in the LSTM at timestep t are:

- Gates:
 - ▶ $\sigma(\mathbb{R}^{(N_{\text{hidden}}, N_{\text{input}} + N_{\text{hidden}})} \cdot \mathbb{R}^{(N_{\text{hidden}} + N_{\text{input}}, 1)}) \rightarrow f_t \in \mathbb{R}^{(N_{\text{hidden}}, 1)}$
 - ▶ Input Gate and Output Gate are similar
- Memory Cell: $C_t \in [(\mathbb{R}^{(N_{\text{hidden}}, 1)} \odot \mathbb{R}^{(N_{\text{hidden}}, 1)}) + (\mathbb{R}^{(N_{\text{hidden}}, 1)} \odot \mathbb{R}^{(N_{\text{hidden}}, 1)})] \rightarrow C_t \in \mathbb{R}^{N_{\text{hidden}}, 1}$

For the backward process, we also use **BPTT algorithm** to train the LSTM.

2. Code Implementation

TBD