

## RNN->LSTM video notes

笔记本： 深度学习

创建时间： 2018/12/17 20:21

更新时间： 2019/7/2 14:23

作者： beyourselfwb@163.com

URL: <https://www.youtube.com/watch?v=WCUNPb-5EYI>

---

<https://www.youtube.com/watch?v=WCUNPb-5EYI>

深度好文：

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

<https://medium.com/mlreview/understanding-lstm-and-its-diagrams-37e2f46f1714>

LSTM VS RNN : solve gradient vanishing/exploding, handling long-term dependencies

LSTM VS CNN : varying-length texts

The fall of RNN / LSTM

来源: <https://towardsdatascience.com/the-fall-of-rnn-lstm-2d1594c74ce0>

LSTM最多也就能记住 100s内的数据, 1000s、10000s搞不定!

LSTM训练很耗时耗资源、也没办法运行在浏览器、手机、智能音响上

Q: How does LSTM help prevent the vanishing (and exploding) gradient problem in a recurrent neural network?

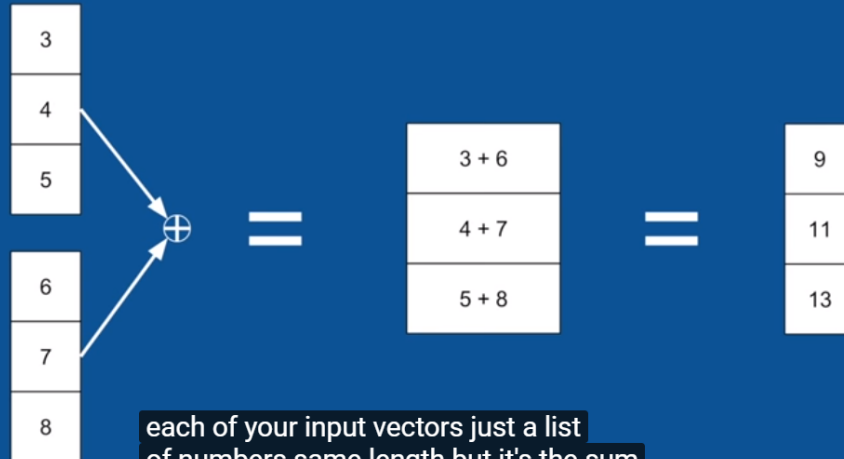
A: Use gating (pass or block, or in other words, 1 or 0) function, not activation function. And train the 'combination' of all those gates. Doing this, no matter how 'deep' your network is, or how 'long' the input sequence is, the network can remember those values, as long as those gates are all 1 along the path. <- This is how LSTM/GRU did the job.

from: <https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network>

## Recurrent Neural Networks (RNN) and Long Short-Term



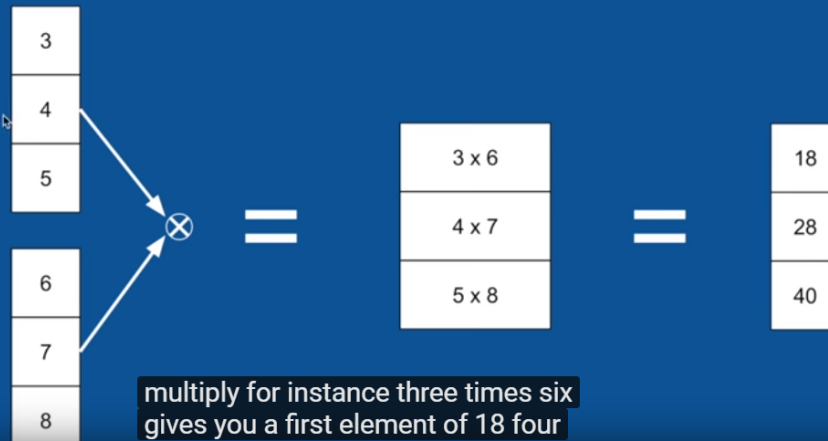
## Plus junction: element-by-element addition



each of your input vectors just a list of numbers same length but it's the sum

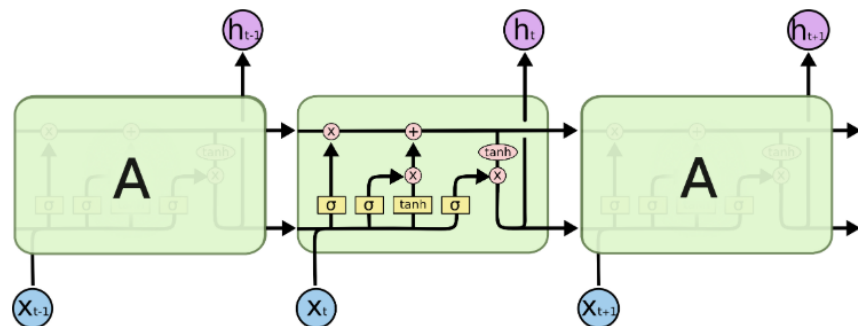
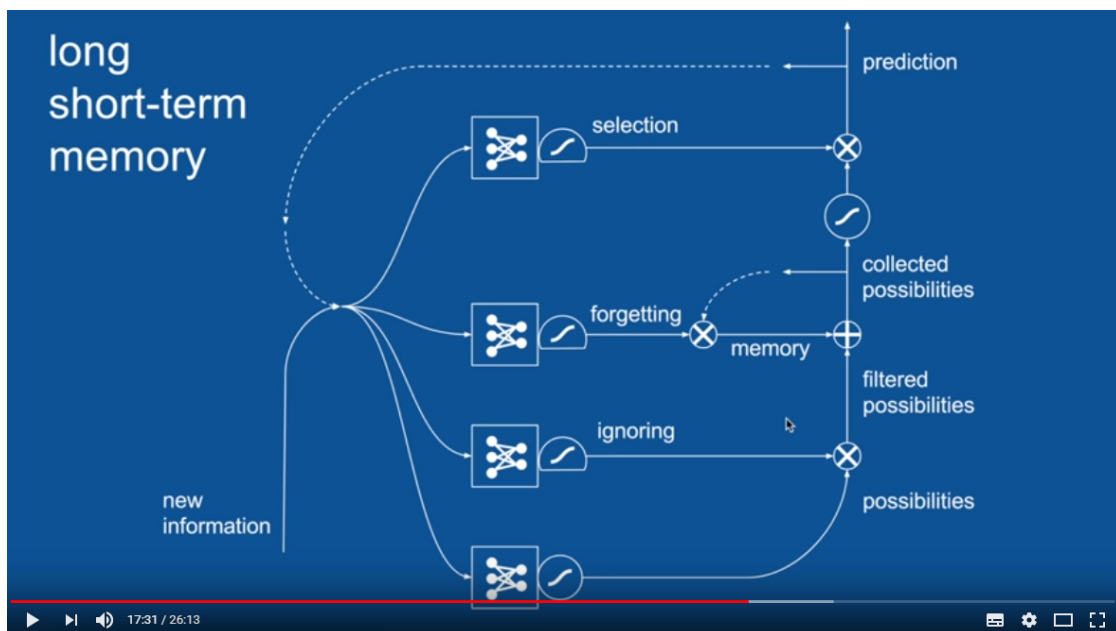
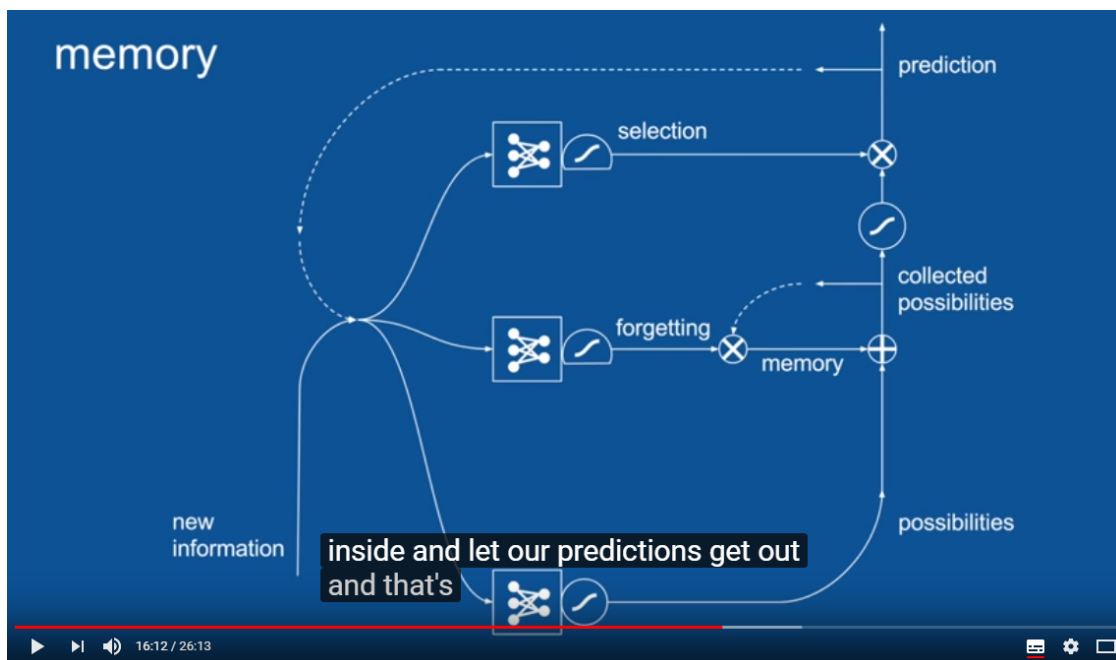
Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM)

## Times junction: element-by-element multiplication



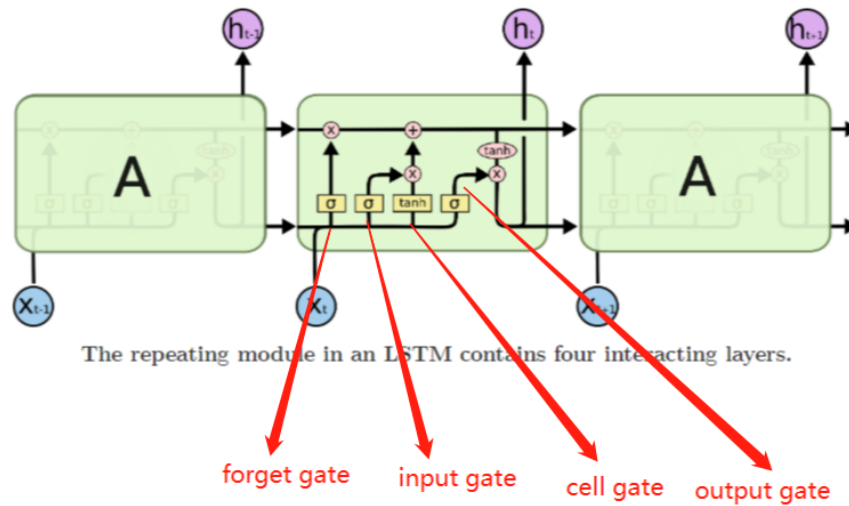
multiply for instance three times six gives you a first element of 18 four

Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM)



The repeating module in an LSTM contains four interacting layers.

放到一起对比



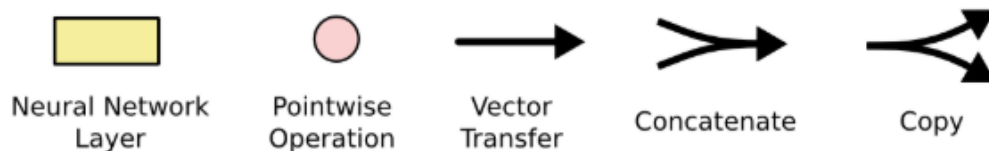
$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$

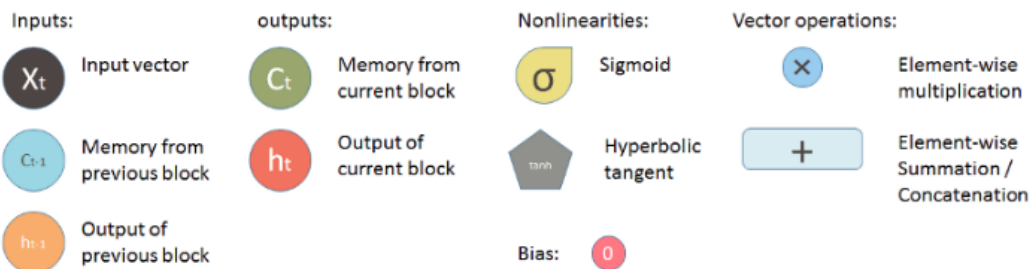
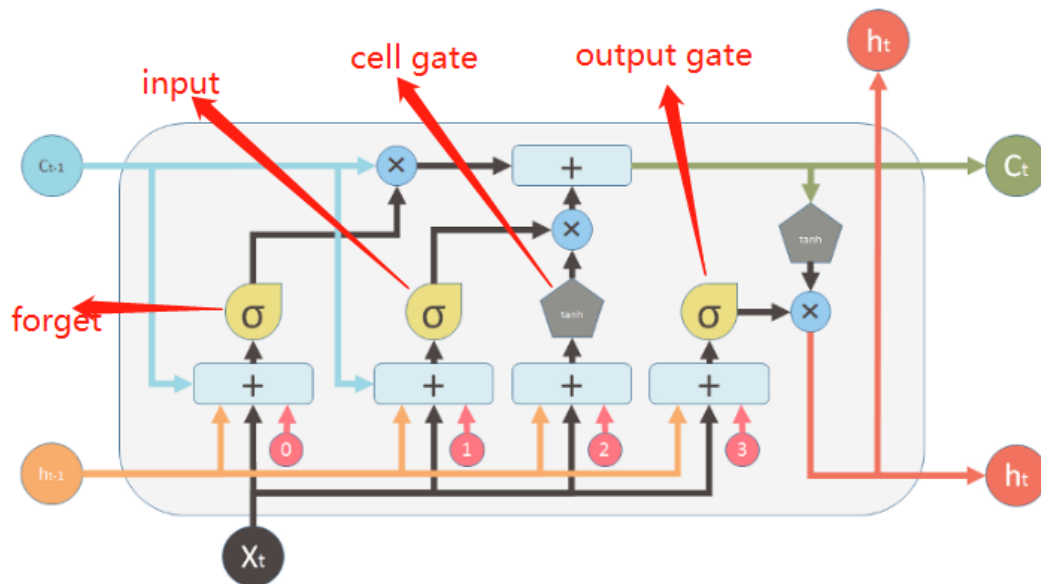
$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$





```
class torch.nn.LSTMCell(input_size, hidden_size, bias=True) \[source\] 🔗
```

A long short-term memory (LSTM) cell.

$$\begin{aligned}
 i &= \text{sigmoid}(W_{ii}x + b_{ii} + W_{hi}h + b_{hi}) \\
 f &= \text{sigmoid}(W_{if}x + b_{if} + W_{hf}h + b_{hf}) \\
 g &= \tanh(W_{ig}x + b_{ig} + W_{hg}h + b_{hg}) \\
 o &= \text{sigmoid}(W_{io}x + b_{io} + W_{ho}h + b_{ho}) \\
 c' &= f * c + i * g \\
 h' &= o * \tanh(c')
 \end{aligned}$$

Parameters:

- **input\_size** – The number of expected features in the input  $x$
- **hidden\_size** – The number of features in the hidden state  $h$
- **bias** – If *False*, then the layer does not use bias weights  $b_{ih}$  and  $b_{hh}$ . Default: `True`

# Sequential patterns

Text

Speech

Audio

Video

Physical processes

Anything embedded in time (almost everything)

those grammar structures that are  
specific to each language and



Traditional LSTM with forget gates.<sup>[2][3]</sup>

Initial values:  $c_0 = 0$  and  $h_0 = 0$ . The operator  $\circ$  denotes the **Hadamard product** (entry-wise product).

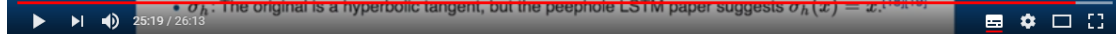
$$\begin{aligned}f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\c_t &= f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \\h_t &= o_t \circ \sigma_h(c_t)\end{aligned}$$

Variables

- $x_t$ : input vector
- $h_t$ : output vector
- $c_t$ : cell state vector
- $W$ ,  $U$  and  $b$ : parameter matrices and vector
- $f_t$ ,  $i_t$  and  $o_t$ : gate vectors
  - $f_t$ : Forget gate vector. Weight of remembering old information.
  - $i_t$ : Input gate vector. Weight of acquiring new information.
  - $o_t$ : Output gate vector. Output candidate.

Activation functions

- $\sigma_g$ : The original is a sigmoid function.
- $\sigma_c$ : The original is a hyperbolic tangent.
- $\sigma_h$ : The original is a hyperbolic tangent, but the peephole LSTM paper suggests  $\sigma_h(x) = x$ .<sup>[1][2]</sup>



lifted straight from the Wikipedia page  
I won't

## Variants [\[edit\]](#)

In the equations below, the lowercase variables represent vectors. Matrices  $W_q$  and  $U_q$  contain, respectively, the weights of the input and recurrent connections, where  $q$  can either be the input gate  $i$ , output gate  $o$ , the forget gate  $f$  or the memory cell  $c$ , depending on the activation being calculated.

### LSTM with a forget gate [\[edit\]](#)

The compact forms of the equations for the forward pass of an LSTM unit with a forget gate are:<sup>[1][2]</sup>

$$\begin{aligned}f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\c_t &= f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \\h_t &= o_t \circ \sigma_h(c_t)\end{aligned}$$

where the initial values are  $c_0 = 0$  and  $h_0 = 0$  and the operator  $\circ$  denotes the **Hadamard product** (element-wise product). The subscript  $t$  indexes the time step.



### LSTM with a forget gate [\[edit\]](#)

The compact forms of the equations for the forward pass of an LSTM unit with a forget gate are:<sup>[1][2]</sup>

$$\begin{aligned}f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\c_t &= f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \\h_t &= o_t \circ \sigma_h(c_t)\end{aligned}$$

where the initial values are  $c_0 = 0$  and  $h_0 = 0$  and the operator  $\circ$  denotes the [Hadamard product](#) (element-wise product). The subscript  $t$  indexes the time step.

#### Variables [\[edit\]](#)

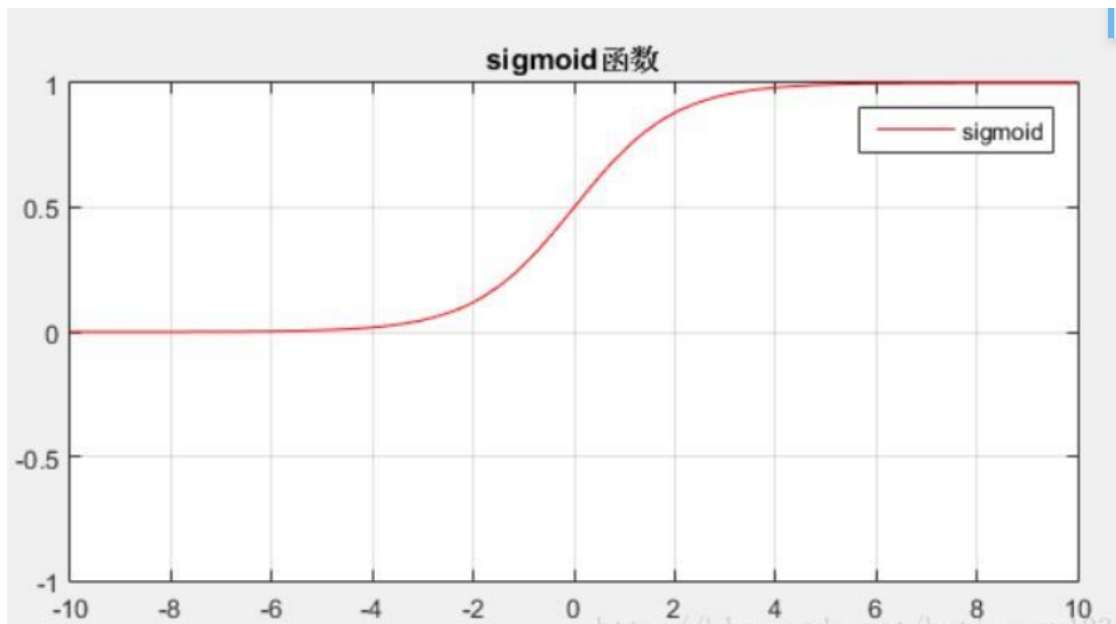
- $x_t \in \mathbb{R}^d$ : input vector to the LSTM unit
- $f_t \in \mathbb{R}^h$ : forget gate's activation vector
- $i_t \in \mathbb{R}^h$ : input gate's activation vector
- $o_t \in \mathbb{R}^h$ : output gate's activation vector
- $h_t \in \mathbb{R}^h$ : hidden state vector also known as output vector of the LSTM unit
- $c_t \in \mathbb{R}^h$ : cell state vector
- $W \in \mathbb{R}^{h \times d}$ ,  $U \in \mathbb{R}^{h \times h}$  and  $b \in \mathbb{R}^h$ : weight matrices and bias vector parameters which need to be learned during training

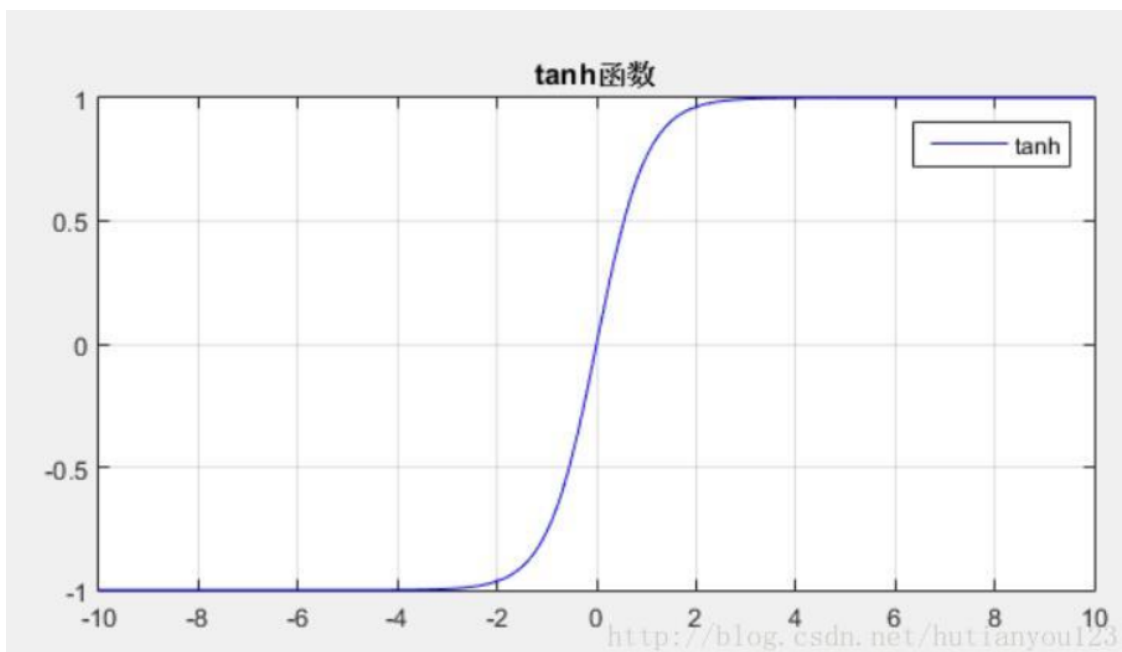
where the superscripts  $d$  and  $h$  refer to the number of input features and number of hidden units, respectively.

#### Activation functions [\[edit\]](#)

- $\sigma_g$ : sigmoid function.
- $\sigma_c$ : hyperbolic tangent function.
- $\sigma_h$ : hyperbolic tangent function or, as the peephole LSTM paper<sup>[which?]</sup> suggests,  $\sigma_h(x) = x$ .<sup>[20][21]</sup>

激活函数：

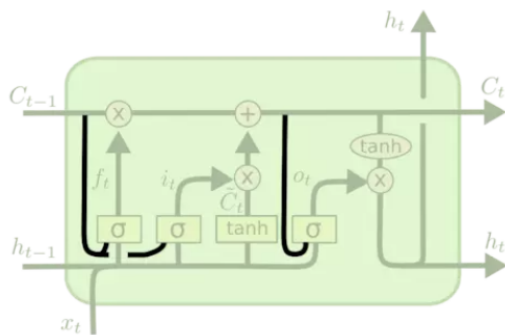




## LSTM 的变体

我们到目前为止都还在介绍正常的 LSTM。但是不是所有的 LSTM 都长成一个样子的。实际上，几乎所有包含 LSTM 的论文都采用了微小的变体。差异非常小，但是也值得拿出来讲一下。

其中一个流形的 LSTM 变体，就是由 [Gers & Schmidhuber \(2000\)](#) 提出的，增加了“**peephole connection**”。是说，我们让门层也会接受细胞状态的输入。



$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

# Resources

Chris Olah's [tutorial](#)

Andrej Karpathy's

[Blog post](#)

[RNN code](#)

[Stanford CS231n lecture](#)

The [DeepLearning 4J](#) tutorial has some helpful discussion and a longer list of good resources.

[How Neural Networks Work](#) [\[video\]](#)