



Eniola Alese

Follow

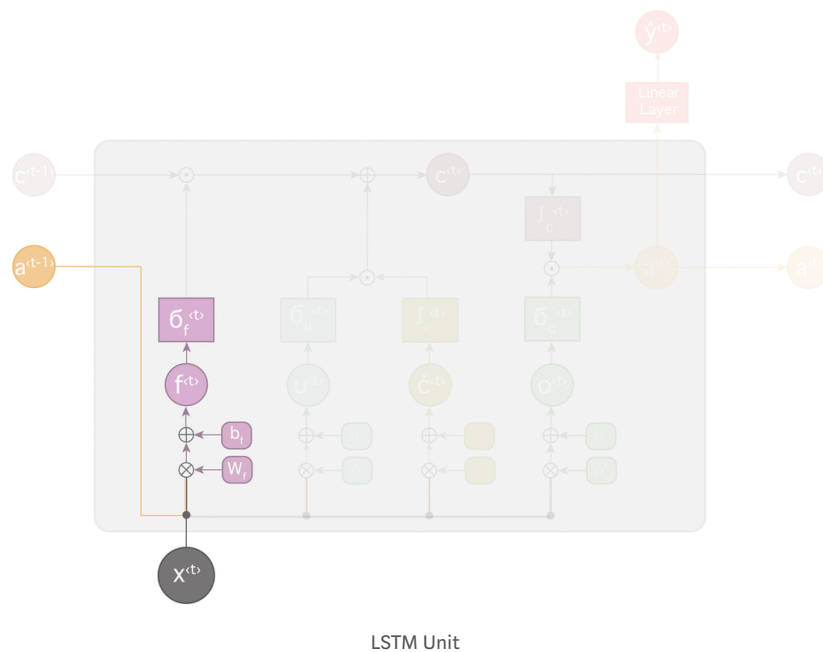
I don't know so I learn and write.

Jun 8 · 4 min read

And of course, LSTM—Part I

The ABC's of the LSTM

In our [last post](#) we looked at the reasons why gradients vanish and presented the LSTM architecture as one of the solutions for dealing with the problem. In this post, we would look at the LSTM's internal structure and gain an intuitive understanding of its computations.



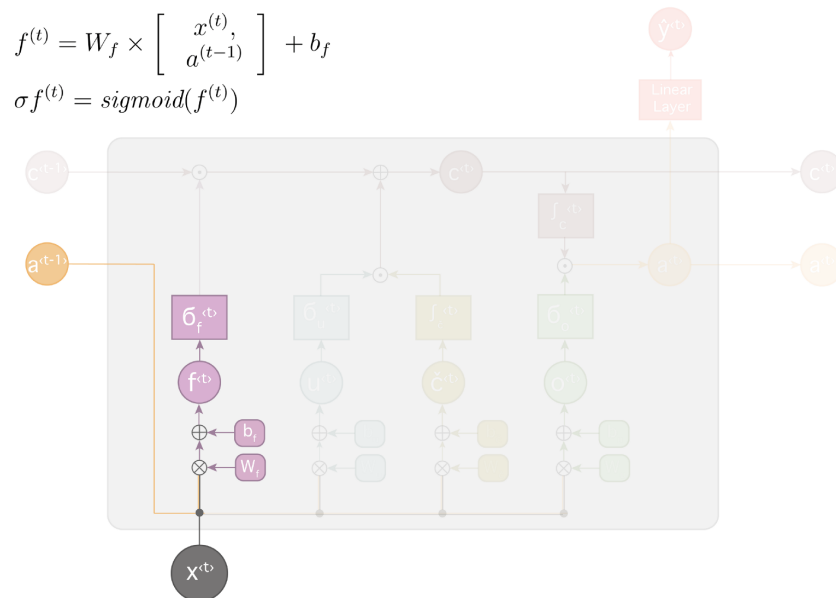
The LSTM, just like in a standard RNN receives its input from the current time-step input $x^{(t)}$ and from the previous time-step hidden state activation $a^{(t-1)}$. The main structural differences between the two units are:

- the introduction of a memory cell state $c^{(t)}$,
- introduction of three sigmoid gates (forget gate $of^{(t)}$, update gate $ou^{(t)}$, output gate $oo^{(t)}$),

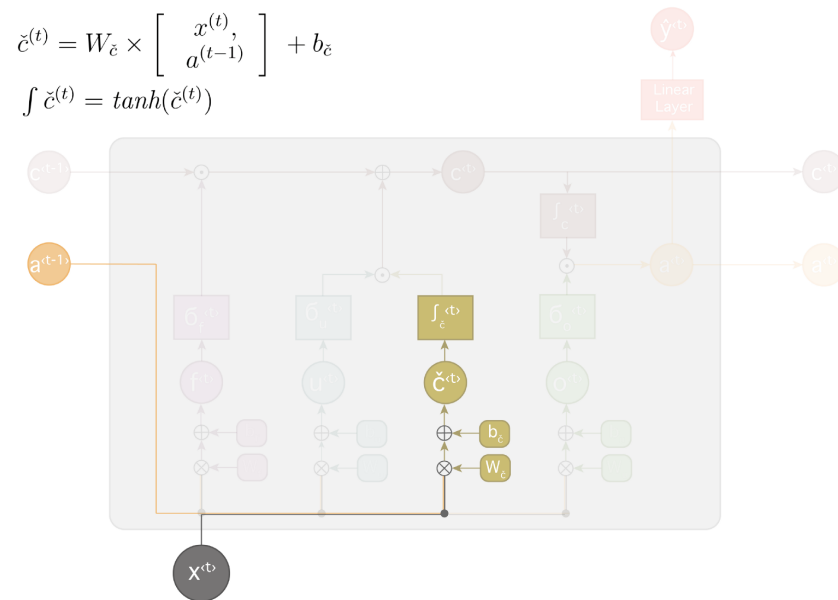
- and the ability to remove or add information to the memory cell state.

Lets break down what happens into individual steps:

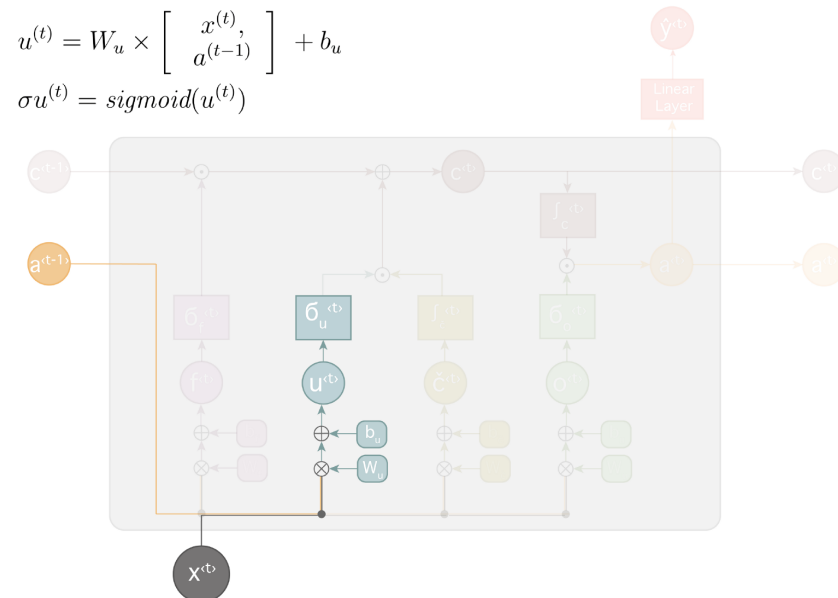
Step 1: The forget gate $\sigma f^{(t)}$ acts like a switch and decides whether or not to discard the memory cell state information that came from the previous time-step $c^{(t-1)}$. It does this by applying a linear transformation to its inputs $x^{(t)}$ and $a^{(t-1)}$, and then passes the output through a logistic function, which then outputs a value between 0 and 1.



Step 2: Next, calculate the temporary memory activation $\tilde{f}^{(t)}$ which holds the temporary memory cell information about the current time-step. This is computed by passing the inputs $x^{(t)}$ and $a^{(t-1)}$ through a linear layer and then through a tanh activation function which then outputs a value between -1 and 1.



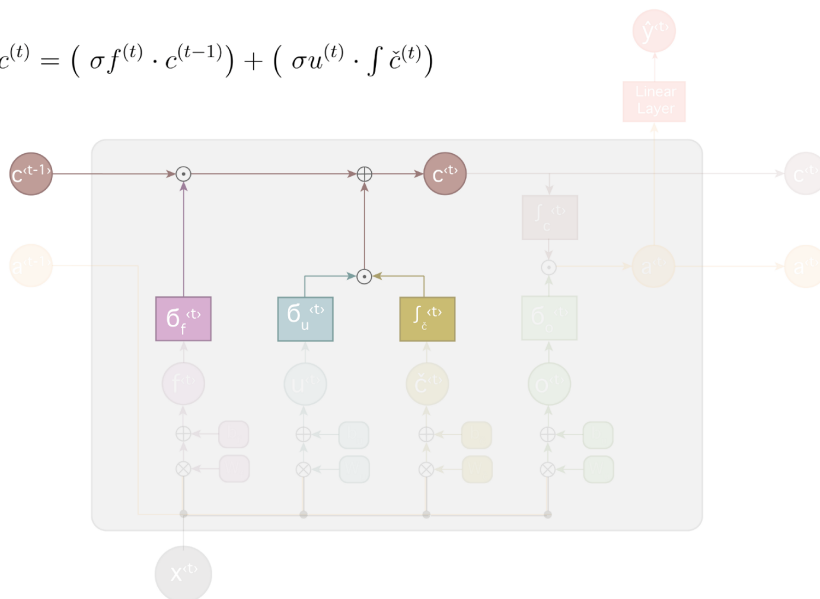
Step 3: Next we have the update gate $\sigma u^{(t)}$ which has the same computational form as the forget gate (i.e linear-to-logistic) and like the forget gate its job is to decide whether or not to discard the temporary memory activation $f_{\tilde{c}}^{(t)}$.



Step 4: Next we update the memory cell state $c^{(t)}$ which is the core of the LSTM design because its what enables the LSTM remember information over long distance. We do this by multiplying the cell state information that came from the previous time-step $c^{(t-1)}$ by the

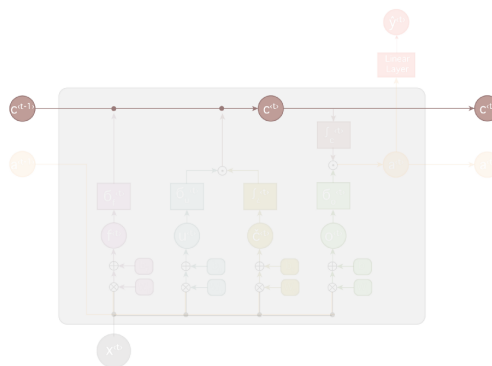
forget gate $\sigma f_{<t>}$ and adding this to the multiplication of the update gate $\sigma u_{<t>}$ and the temporary memory activation $\int \tilde{c}_{<t>}$.

$$c^{(t)} = (\sigma f^{(t)} \cdot c^{(t-1)}) + (\sigma u^{(t)} \cdot \int \tilde{c}^{(t)})$$

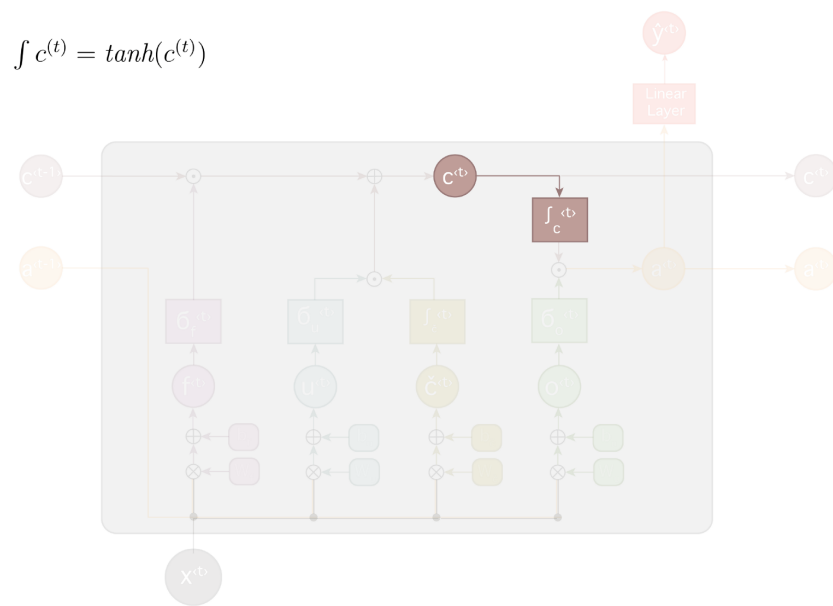


Depending on the values of the forget and update gates the memory cell state can exhibit various behaviours, like if the forget gate $\sigma f_{<t>}$ is ON (1) and the update gate $\sigma u_{<t>}$ is OFF (0) the memory cell state would compute something close to the identity function: $c_{<t>} = c_{<t-1>}$

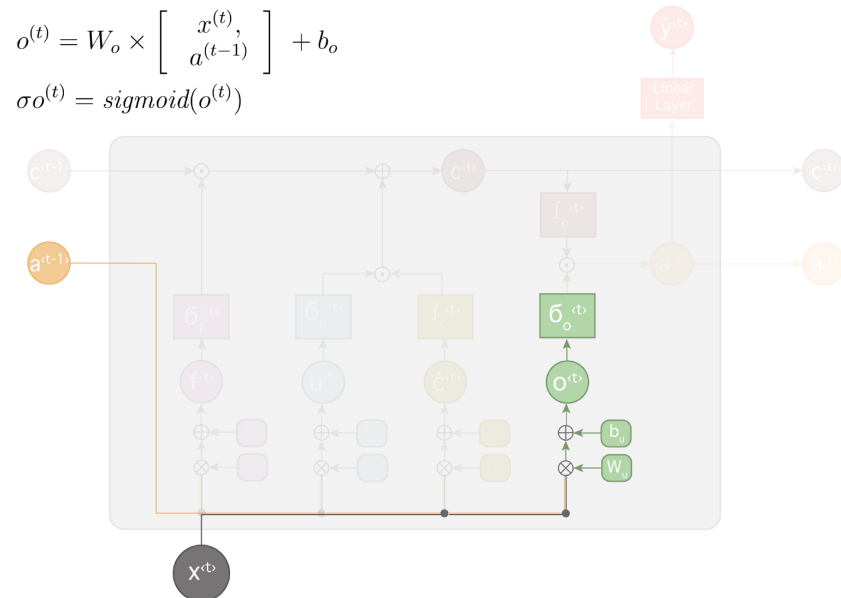
$\sigma f^{(t)}$	$\sigma u^{(t)}$	$c^{(t)}$
0	0	0
1	0	$c^{(t-1)}$
0	1	$\int \tilde{c}^{(t)}$
1	1	$c^{(t-1)} + \int \tilde{c}^{(t)}$



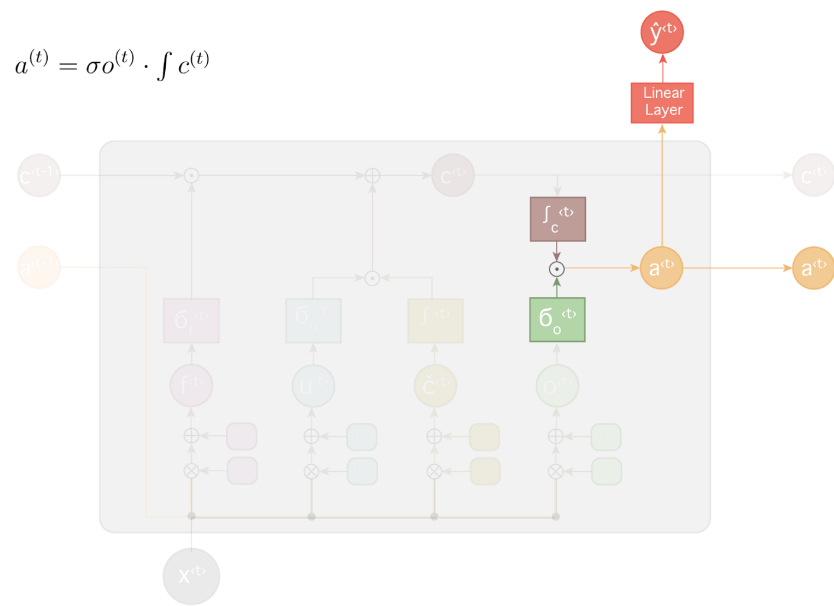
Step 5: Next, we compute the memory activation $\int c_{<t>}$ which is simply the memory cell state $c_{<t>}$ passed through a tanh function.



Step 6: Next, we have the output gate $\sigma_o<t>$ which like the forget and update gates has the same linear-to-logistic computational form, and determines whether or not the memory activation $f_c<t>$ would be passed to the rest of the network.



Step 7: Finally, we compute the hidden state activation $a<t>$ by multiplying the output gate $\sigma_o<t>$ with the memory activation $f_c<t>$ and pass on the output to the next time step and to the predicted output linear layer.



Conclusion

The steps above basically cover all that happens during forward propagation in a single LSTM unit. In the second part of this post we would look at deriving the LSTM gradients for back propagation and a practical implementation of the model in PyTorch.

