

THE NEURAL PERSPECTIVE

DEEP LEARNING SIMPLIFIED

HOME · TUTORIALS · READINGS · CONTACT · ABOUT

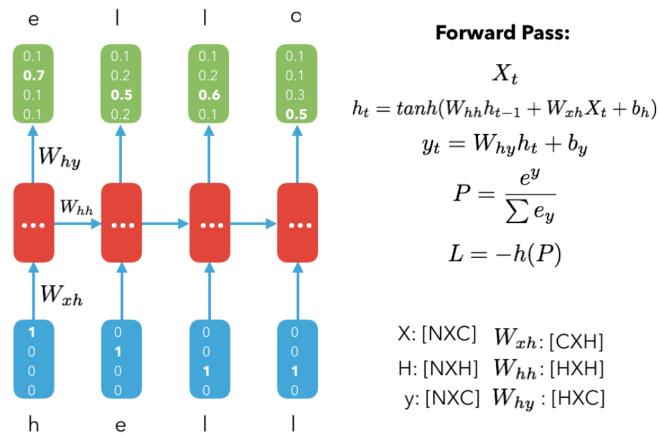
RECURRENT NEURAL NETWORK (RNN) – PART 5: CUSTOM CELLS

Edit

In this post, we will explore the idea of creating our own custom RNN cells. But first, we will take a closer look at the simple RNN and then more complicated units such as LSTM and GRU. We will also analyze the tensorflow code for these units and draw from them to eventually create our own custom cells. In this post, I will be using images from one of the best posts out there on RNNs/LSTMs by Chris Olah. I highly urge you to read the [post](#) and in my post I will be reiterating a lot of material but I will move rather quickly and focus more on the tf code. I will be referring back to this code in a future post on applying layer normalization to these RNN architectures, which can be found [here](#).

BASIC RNNs:

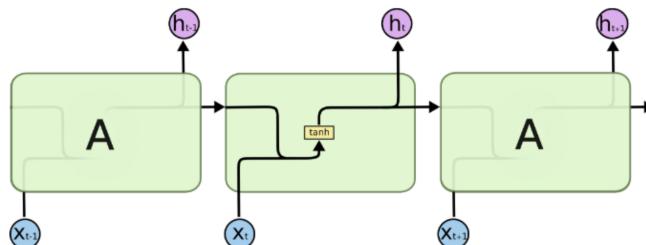
With traditional RNNs, the main issue is that we cannot adequately learn long term dependencies because the operations that we repeat at each cell unit for each input are static. If you think back to the basic RNN cell, the operations all involve the single tanh operation.



This architecture is suitable for inputs where the solutions are based on short term dependencies but if we wish to utilize long term memory efficiently to predict the right targets, we will need a RNN cell unit that is more robust. Cue the LSTM.

LONG SHORT TERM MEMORY NETWORKS (LSTMS):

The architecture of the LSTMs allows us to have long term information control at the expense of more operations. Our traditional RNNs had one output which served as both the hidden state representation and the output from the cell.



There is an absence of information control with this basic architecture that prevents us from holding on to useful information for many steps down the line. The LSTM, instead, has two different types of outputs. We still have the traditional state output which acts as the hidden state representation and the cell's output but the cell also outputs a cell state C . Here is the LSTM in all its glory, time to break it down into pieces.

- SEARCH -

Search ...

- FOLLOW ME ON TWITTER -

Tweets by [@GokuMohandas](#)

 **Goku Mohandas**
@GokuMohandas

Your model may be performing really well by incorrectly focusing on confounding features (extraneous influencers in the data that aren't accounted for). Check out [@johnrzech](#)'s post where x-ray stickers unintentionally influenced the classifications: [medium.com/@jrzech/what-a...](#)



Embed View on Twitter

- RECENT POSTS -

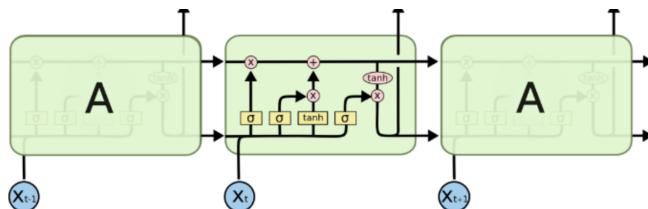
[Update on Embeddings \(Spring 2017\)](#)

[Exploring Sparsity in Recurrent Neural Networks](#)

[Question Answering from Unstructured Text by Retrieval and Comprehension](#)

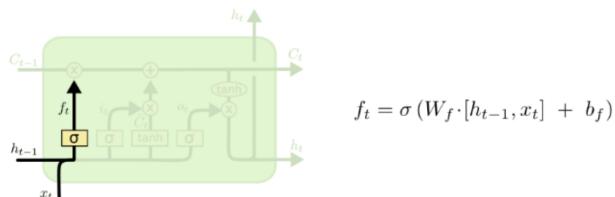
[Overcoming Catastrophic Forgetting in Neural Networks](#)

[Opening the Black Box of Deep Neural Networks via Information](#)



FORGET GATE:

The very first gate is the forget gate. This gate allows us to selectively pass information to determination of the cell state. I will break down the notation below once and you can reapply for all the other gates as well.



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

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

$$W_f : [D+H \times H]$$

$$h_{t-1} : [N \times H]$$

$$x_t : [N \times D]$$

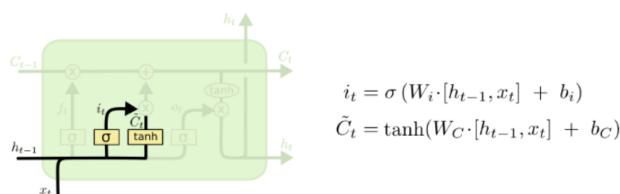
$$b_f : [H]$$

And of course to implement this, you could follow something like tf's `_linear` function. But the main idea is that we are applying this sigmoid operation to both the input and the previous hidden state. But what exactly is applying this sigmoid operation doing? Recall that sigmoid outputs in the range $[0, 1]$ and here we are applying it to a matrix of shape $[N \times H]$, which means we will produce $N \times H$ values with sigmoid applied to them. If the sigmoid operation results in 0, then that hidden value is nullified and if it is 1, we completely let that value be used. Anything in between allows parts of the information to go through. This is a nice way to control the information that is flowing through by effectively blocking and selectively passing parts of the inputs to the cell.

This forget gate, however, is only the first operation that we do to ultimately calculate our cell state. The next operation involves the input gate.

INPUT GATE:

The input gate takes in our input X and the previous hidden state and computes two operations. First it selectively allows parts of the inputs to pass through with a sigmoid gate and then we multiply it by the tanh of the inputs.

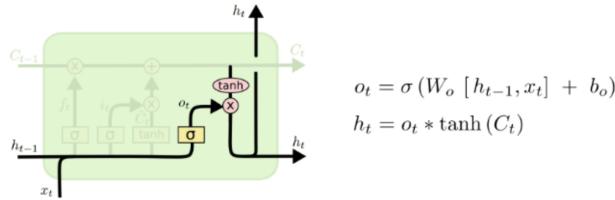


What the tanh is doing here is a bit different from the sigmoid operation. Recall that tanh changes our inputs into the range $[-1, 1]$. This essentially changes the underlying representation of our inputs with this nonlinearity. This is the exact same step as what we were doing with the basic RNN cell. But now we take the product of these two values and add it to the value from the forget gate to calculate our cell state.

These operations with the forget and input gate can be translated to the fact that we keep parts of the old cell state (C_{t-1}) and keep parts of the new transformed (tanh) cell state \tilde{C}_t . These weights are trained with our data to learn exactly how much information to keep and how to perform the correct transformation.

OUTPUT GATE:

The last gate is the output gate and it uses the input, previous hidden state and the new cell state to determine the new hidden state representation.



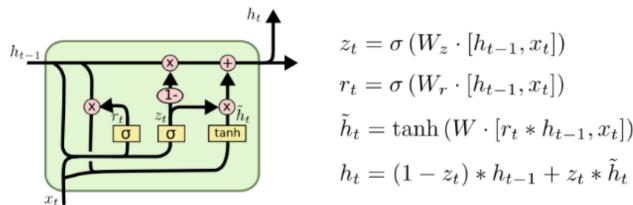
This operation again involves the selective information barrier sigmoid which is multiplied with tanh of the cell state. Note that this tanh operation is not a neural network as with the tanh operation in the input gate. This is simply applying the tangent to the cell state without any modifications with weights. We are merely forcing the cell states [NXH] values to be in the range [-1, 1].

VARIATIONS:

There are literally hundreds of variations for RNN cells so I suggest checking our Chris Olah's [blog](#) again for more information. A few note worthy one's he discussed were the peephole model (allow all gates to see the cell state available at that point in time (C_{t-1} or C_t is already calculated) and coupled cell states (only update when we forget and forget when we update). But the current rival to the LSTM, which is heavily based off of the LSTM and it rapidly growing in use is the Gated Recurrent Unit (GRU).

GATED RECURRENT UNIT (GRU):

The main idea behind the GRU is that it combines the forget and input gate into one update gate.



Empirically, the GRU's performance on most tasks is on par with the LSTM and also computationally less expensive. These tradeoffs are the reason behind its surging popularity.

TENSORFLOW NATIVE IMPLEMENTATIONS:

Now we will take a look at the official Tensorflow code the GRU unit and we will mostly focus on the function calls, inputs and outputs. From here, we will replicate the structure to create our own unique cells. If you're interested in the other cells available, you can find them all at this [link](#). We will just focus on the GRU because its performance is as good as the LSTM in more cases and significantly less complex.

```

1  class GRUCell(RNNCell):
2      """Gated Recurrent Unit cell (cf. http://arxiv.org/abs/1406.1078)"""
3
4      def __init__(self, num_units, input_size=None, activation=None):
5          if input_size is not None:
6              logging.warn("%s: The input_size parameter is deprecated.", self)
7              self._num_units = num_units
8              self._activation = activation
9
10     @property
11     def state_size(self):
12         return self._num_units
13
14     @property
15     def output_size(self):
16         return self._num_units
17
18     def __call__(self, inputs, state, scope=None):
19         """Gated recurrent unit (GRU) with nunits cells."""
20         with vs.variable_scope(scope or type(self).__name__):
21             with vs.variable_scope("Gates"):
22                 # We start with bias of 1.0 to not reset and not update.
23                 r, u = array_ops.split(1, 2, _linear([inputs, state],
24                                                 2 * self._num_units))
25                 r, u = sigmoid(r), sigmoid(u)
26             with vs.variable_scope("Candidate"):
27                 c = self._activation(_linear([inputs, r * state],
28                                              self._num_units, True))
29                 new_h = u * state + (1 - u) * c
30
31         return new_h, new_h

```

The GRUCell class starts with the `__init__` function which defines the number of units and the activation function it will use. This is the activation function that is usually tanh but the sigmoid activations are fixed since the [0,1] range allows us to control the information flow. Then we have two properties that both return `self._num_units` when invoked. And finally, we have our `__call__` function which is what processes the input and churns out the new hidden state. Recall that GRU does not have a cell state like the LSTM.

First, we compute r and u ($u = z$ in colah's notation above). Instead of separately doing them, we just merge the weights and do it with $2 * \text{num_units}$ and then we split it by two. **split(dim, num_splits, value)**. Then we apply our sigmoid activate on the values to selectively control the information flow. Then we calculate the candidate c and use it to calculate out new hidden state representation. You may see that the order for calculating new_h is switched, either way works fine, because the weights will train accordingly.

All of the other cells' codes look very similar to this, so you will easily be able to interpret them.



Posted in: Machine Translation, Optimization/Architecture, Sequence-to-Sequence

Tagged: Tutorials

← Embeddings (skipgram and
CBOW) Implementations

RECURRENT NEURAL NETWORKS (RNN) – PART
3: Encoder-Decoder →

ONE THOUGHT ON “RECURRENT NEURAL NETWORK (RNN) – PART 5: CUSTOM CELLS”

Pingback: RECURRENT NEURAL NETWORKS (RNN) – PART 3: Neural Machine Translation (NMT) –
The Neural Perspective Edit

LEAVE A REPLY

Enter your comment here...

