# Building a Long-Short Term Memory (LSTM) Unit

**Vaishnavi Shrivastava**

**Github Link:** https://github.com/vshrivas/CogNETive

# Forward Propagation

Forward propagation is defined in the *forwardprop* method, which takes in $h_{t-1}$ as $h_{prev}$, $C_{t-1}$ as $c_{prev}$, and the current input $x$.
$z = [h_{t-1}, x_t]$, concatenation of $h_{t-1}$ and $x_t$.

## Forget Gate

The forget gate layer uses $h_{t-1}$ and $x_t$ to calculate a number between 0 and 1 for each value in $C_{t-1}$, where 0 means to completely forget and 1 means to completely remember. Sigmoid is used because it outputs a value between 0 and 1, which corresponds to an amount of information flow.

$$f_t = \sigma(W_f \cdot z + b_f)$$

## Input Gate and Cell State

This is used to decide the values that need to be updated, and by how much they need to be updated. We also generate a set of candidate values $\tilde{C}_t$ that can be multiplied by the input gate and added to the new cell state.
$i_t = \sigma(W_i \cdot z + b_i)$

The cell state's values should be able to increase or decrease when we add the output of some function. Sigmoid output is always non-negative; meaning values in the state would only increase. The output from tanh can be positive or negative, allowing for increases and decreases in the state. That's why tanh is used to determine candidate values to get added to the internal state. On the other hand, to overcome the vanishing gradient problem, we need a function whose second derivative can sustain for a long range before going to zero. Tanh is a good function with the above property.
A tanh function ensures that the values stay between -1 and 1, thus regulating the output of the neural network.
$\tilde{C}_t = \tanh{(W_C \cdot z + b_c)}$

We multiply the old state by $f_t$, forgetting the things we decided to forget earlier. Then we add $i_t \cdot \tilde{C}_t$, the new candidate values, scaled by how much we decided to update each state value.
$C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t$.

*Sidenote:*
In RNNs the weights are shared across the time steps. If you have a ReLU activation and the recurrence matrices enter the unstable regime (for Elman RNNs this occurs if the spectral norm of the rec. matrix is ¿ 1, for LSTMs the condition is more complicated, but still exists), the magnitude of the state will increase exponentially with the time steps during the forward pass. In turn this will increase the the magnitude of the gradients during the backward pass, the two effects can easily feed into each other, leading to numeric

overflow. Properly placed tanh functions can avoid this: at worst you will get saturation and vanishing gradients, which is better than overflow. There are some papers that used ReLUs with RNNs, but they require very careful weight initialization and low learning rates, and in general there is no benefit over LSTMs or GRUs.

## Output Gate and Hidden State

The output from this timestep, $h_t$, will be a filtered version of the cell state, where not all of the information passes through. To decide the information that will pass through, we calculate a sigmoid layer $o_t$. The cell state is pushed through a tanh, to squeeze it between 1 and -1, and then multiplied by $o_t$ to output only the parts we decided.

$o_t = \sigma(W_o \cdot z + b_o)$

$h_t = o_t \cdot tanh(C_t)$

These values are all returned in the dictionary *fresults*.

# Backpropagation Through Time Derivation

`dv =`

$\frac{dL}{dv_t} =$

$\hat{y}_t - y_t$

`dh =`

$\frac{dL}{dh_t} =$

$\frac{dL}{dv_t} \cdot \frac{dv_t}{dh_t} + \frac{dL_{t+1}}{dh_t} =$

$\frac{dL}{\cancel{dv_t}} \cdot \frac{\cancel{dv_t}}{dh_t} + \frac{dL_{t+1}}{dh_t} =$

$\frac{dL}{dv_t} \cdot W_v + \frac{dL_{t+1}}{dh_t} =$

`dv` $\cdot W_v +$ `dh`$_{\text{prev}}$

The second term is added because $h_{t-1}$ is passed to the cell at the next timestep, and therefore has an effect on the cell in future timesteps that needs to be accounted for.

`dc =`

$\frac{dL}{dC_t} =$

$\frac{dL}{dh_t} \cdot \frac{dh_t}{dC_t} + \frac{dL}{dC_{t+1}} \cdot \frac{dC_{t+1}}{dC} =$

$\frac{dL}{\cancel{dh_t}} \cdot \frac{\cancel{dh_t}}{dC_t} + \frac{dL}{dC_{t+1}} \cdot \frac{\cancel{dC_{t+1}}}{dC} =$

$\frac{dL}{dh_t} \cdot [(o_t)(\tanh' C_t)] + \frac{dL}{dC_{t+1}} \cdot f_{t+1} =$

`dh` $\cdot [o_t(1 - \tanh^2 C_t)] +$ `dc`$_{\text{prev}}$

$C_t$ impacts $L$ by influencing $h_t$ in the current timestep, and by influencing $C_{t+1}$ in the next timestep, so we add both factors to find the derivative.

## Forget Gate

`df =`

$\frac{dL}{df_t} =$

$\frac{dL}{dC_t} \cdot \frac{dC_t}{df_t} =$

$\frac{dL}{\cancel{dC_t}} \cdot \frac{\cancel{dC_t}}{df_t} =$

`dc` $\cdot C_{t-1}$

The `logit` function is the inverse of the `sigmoid` function, so applying it to $f_t$ will yield $W_f Z + b_f$. Let

`dsigmoid` be the derivative of the sigmoid function.

$\texttt{dW}_\texttt{f} =$

$\frac{dL}{dW_f} =$

$\frac{dL}{df_t} \cdot \frac{df_t}{dW_f} =$

$\frac{dL}{d\!\!\!/f_t} \cdot \frac{d\!\!\!/f_t}{dW_f} =$

$\frac{dL}{df_t} \cdot [\sigma'(W_f Z + b_f)(Z)] =$

$\texttt{df} \cdot [\texttt{dsigmoid}(\texttt{logit}(f))] \cdot Z$

$\texttt{db}_\texttt{f} =$

$\frac{dL}{db_f} =$

$\frac{dL}{df_t} \cdot \frac{df_t}{db_f} =$

$\frac{dL}{d\!\!\!/f_t} \cdot \frac{d\!\!\!/f_t}{db_f} =$

$\frac{dL}{df_t} \cdot [\sigma'(W_f Z + b_f)(1)] =$

$\texttt{df} \cdot [\texttt{dsigmoid}(\texttt{logit}(f))]$

$\texttt{dZ}_\texttt{f} =$

$\frac{dL}{dZ} =$

$\frac{dL}{df_t} \cdot \frac{df_t}{dZ} =$

$\frac{dL}{d\!\!\!/f_t} \cdot \frac{d\!\!\!/f_t}{dZ} =$

$\frac{dL}{df_t} \cdot [\sigma'(W_f Z + b_f)(W_f)] =$

$\texttt{df} \cdot [\texttt{dsigmoid}(\texttt{logit}(f))] \cdot W_f$

In the code, the repeated multiplication of `df` and $[\texttt{dsigmoid}(\texttt{logit}(f))]$ is eliminated by setting `df = df` $\cdot [\texttt{dsigmoid}(\texttt{logit}(f))]$.

## Input Gate

The `logit` function is the inverse of the `sigmoid` function, so applying it to $i_t$ will yield $W_i Z + b_i$. Let `dsigmoid` be the derivative of the sigmoid function.

$\texttt{di} =$

$\frac{dL}{di_t} =$

$\frac{dL}{dC_t} \cdot \frac{dC_t}{di_t} =$

$\frac{dL}{d\tilde{C}_t} \cdot \frac{d\tilde{C}_t}{di_t} =$

$\texttt{dc} \cdot \tilde{C}_t$

$\texttt{dW}_\texttt{i} =$

$\frac{dL}{dW_i} =$

$\frac{dL}{di_t} \cdot \frac{di_t}{dW_i} =$

$\frac{dL}{d\!\!\!/i_t} \cdot \frac{d\!\!\!/i_t}{dW_i} =$

$\frac{dL}{di_t} \cdot [\sigma'(W_i Z + b_i)(Z)] =$

$\texttt{di} \cdot [\texttt{dsigmoid}(\texttt{logit}(i_t))] \cdot Z$

$\texttt{db}_\texttt{i} =$

$\frac{dL}{db_i} =$

$\frac{dL}{di_t} \cdot \frac{di_t}{db_i} =$

$\frac{dL}{d\!\!\!/i_t} \cdot \frac{d\!\!\!/i_t}{db_f} =$

3

$\frac{dL}{di_t} \cdot [\sigma'(W_i Z + b_i)(1)] =$
di $\cdot [\text{dsigmoid}(\text{logit}(i_t))]$

dZ$_i$ =
$\frac{dL}{dZ} =$
$\frac{dL}{di_t} \cdot \frac{di_t}{dZ} =$
$\frac{dL}{d\not{i_t}} \cdot \frac{d\not{i_t}}{dZ} =$
$\frac{dL}{di_t} \cdot [\sigma'(W_i Z + b_i)(W_i)] =$
df $\cdot [\text{dsigmoid}(\text{logit}(i_t))] \cdot W_i$

In the code, the repeated multiplication of di and $[\text{dsigmoid}(\text{logit}(i_t))]$ is eliminated by setting
di = di $\cdot [\text{dsigmoid}(\text{logit}(i_t))]$.

## Output Gate

The logit function is the inverse of the sigmoid function, so applying it to $o_t$ will yield $W_o Z + b_o$. Let
dsigmoid be the derivative of the sigmoid function.
do =
$\frac{dL}{do_t} =$
$\frac{dL}{dh_t} \cdot \frac{dh_t}{do_t} =$
$\frac{dL}{d\not{h_t}} \cdot \frac{d\not{h_t}}{do_t} =$
dh $\cdot \tanh C_t$

dW$_o$ =
$\frac{dL}{dW_o} =$
$\frac{dL}{do_t} \cdot \frac{do_t}{dW_o} =$
$\frac{dL}{d\not{o_t}} \cdot \frac{d\not{o_t}}{dW_o} =$
$\frac{dL}{do_t} \cdot [\sigma'(W_o Z + b_o)(Z)] =$
do $\cdot [\text{dsigmoid}(\text{logit}(o_t))] \cdot Z$

db$_o$ =
$\frac{dL}{db_o} =$
$\frac{dL}{do_t} \cdot \frac{do_t}{db_o} =$
$\frac{dL}{d\not{o_t}} \cdot \frac{d\not{o_t}}{db_o} =$
$\frac{dL}{do_t} \cdot [\sigma'(W_o Z + b_o)(1)] =$
do $\cdot [\text{dsigmoid}(\text{logit}(o_t))]$

dZ$_o$ =
$\frac{dL}{dZ} =$
$\frac{dL}{do_t} \cdot \frac{do_t}{dZ} =$
$\frac{dL}{d\not{o_t}} \cdot \frac{d\not{o_t}}{dZ} =$
$\frac{dL}{do_t} \cdot [\sigma'(W_o Z + b_o)(W_o)] =$
do $\cdot [\text{dsigmoid}(\text{logit}(o_t))] \cdot W_o$

In the code, the repeated multiplication of do and $[\text{dsigmoid}(\text{logit}(o_t))]$ is eliminated by setting
do = do $\cdot [\text{dsigmoid}(\text{logit}(o_t))]$.

## Candidate Gate

The `arctanh` function is the inverse of the `tanh` function, so applying it to $c_{bar}$ will yield $W_{cbar}Z + b_{cbar}$. Let `dtanh` be the derivative of the tanh function.

$$\texttt{dc}_{\texttt{bar}} =$$
$$\frac{dL}{dC_{bar_t}} =$$
$$\frac{dL}{dC_t} \cdot \frac{dC_t}{dC_{bar_t}} =$$
$$\frac{dL}{d\cancel{C_t}} \cdot \frac{d\cancel{C_t}}{dC_{bar_t}} =$$
$$\texttt{dc} \cdot i_t$$

$$\texttt{dW}_{\texttt{cbar}} =$$
$$\frac{dL}{dW_{cbar}} =$$
$$\frac{dL}{dC_{bar_t}} \cdot \frac{dC_{bar_t}}{dW_{cbar}} =$$
$$\frac{dL}{d\cancel{C_{bar_t}}} \cdot \frac{d\cancel{C_{bar_t}}}{dW_{cbar}} =$$
$$\frac{dL}{dC_{bar_t}} \cdot [\tanh'(W_{cbar}Z + b_{cbar})(Z)] =$$
$$\texttt{dc}_{\texttt{bar}} \cdot [\texttt{dtanh}(\texttt{arctanh}(c_{bar}))] \cdot Z$$

$$\texttt{db}_{\texttt{cbar}} =$$
$$\frac{dL}{db_{cbar}} =$$
$$\frac{dL}{dC_{bar_t}} \cdot \frac{dC_{bar_t}}{db_{cbar}} =$$
$$\frac{dL}{d\cancel{C_{bar_t}}} \cdot \frac{d\cancel{C_{bar_t}}}{db_{cbar}} =$$
$$\frac{dL}{dC_{bar_t}} \cdot [\tanh'(W_{cbar}Z + b_{cbar})(1)] =$$
$$\texttt{dc}_{\texttt{bar}} \cdot [\texttt{dtanh}(\texttt{arctanh}(c_{bar}))]$$

$$\texttt{dZ}_{\texttt{cbar}} =$$
$$\frac{dL}{dZ} =$$
$$\frac{dL}{dC_{bar_t}} \cdot \frac{dC_{bar_t}}{dZ} =$$
$$\frac{dL}{d\cancel{C_{bar_t}}} \cdot \frac{d\cancel{C_{bar_t}}}{dZ} =$$
$$\frac{dL}{dC_{bar_t}} \cdot [\tanh'(W_{cbar}Z + b_{cbar})(W_{cbar})] =$$
$$\texttt{dc}_{\texttt{bar}} \cdot [\texttt{dtanh}(\texttt{arctanh}(c_{bar}))] \cdot W_{cbar}$$

In the code, the repeated multiplication of $\texttt{dc}_{\texttt{bar}}$ and $[\texttt{dtanh}(\texttt{arctanh}(c_{bar}))]$ is eliminated by setting $\texttt{dc}_{\texttt{bar}} = \texttt{dc}_{\texttt{bar}} \cdot [\texttt{dtanh}(\texttt{arctanh}(c_{bar}))]$.

## Training Process

The network is trained in steps, where in each step, first we run forward propagation through all of the training examples, and store the results of forward prop at each timestep in a dictionary and accumulate the overall loss by comparing the output of each timestep, with the expected result at each timestep. The we perform backpropagation through each timestep, starting from the last timestep and going backwards to the first, calculating and accumulating the gradients of the weight matrices and biases at each step, and sending back the calculated values of $\texttt{dh}_{\texttt{next}}$ and $\texttt{dc}_{\texttt{next}}$ to the previous timestep to be used in the gradient calculations there, since these are the values that link backprop across timesteps.