

GRUS

Each GRU unit, has a C =memory cell, which provides a bit of memory to be later checked:

$$C^{<+>} = a^{<+>}$$

At every timestep, we consider overwriting the memory cell with a value $\hat{C}^{<+>}$:

$$\hat{C}^{<+>} = \tanh(W_c[c^{<t-1>}, x^{<+>}] + b_c)$$

The important idea is the GATE:

- Update gate $\Rightarrow \Gamma_u = o(W_u[c^{<t-1>}, x^{<+>}] + b_u)$

It finally decides whether or not we finally update $\hat{C}^{<+>}$.

So, for example, let's say that $c^{<+>} = 1$ for cat, as it is singular. Then, was is used because it is stored in the $c^{<+>}$ cell.

$c^{<+>} = 1 \dots \dots \dots = 1$
The cat, which already ate ..., was full.

But, $c^{<+>}$ may be required to change its value at another time. That is what Γ_u decides.

Finally, we calculate:

$$C^{<t>} = \Gamma_u \cdot \hat{C}^{<t>} + (1 - \Gamma_u) \cdot C^{<t-1>}$$

Full GRU

How relevant is $C^{<t-1>}$ to computing the next candidate $C^{<t>}?$

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$
$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$
$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

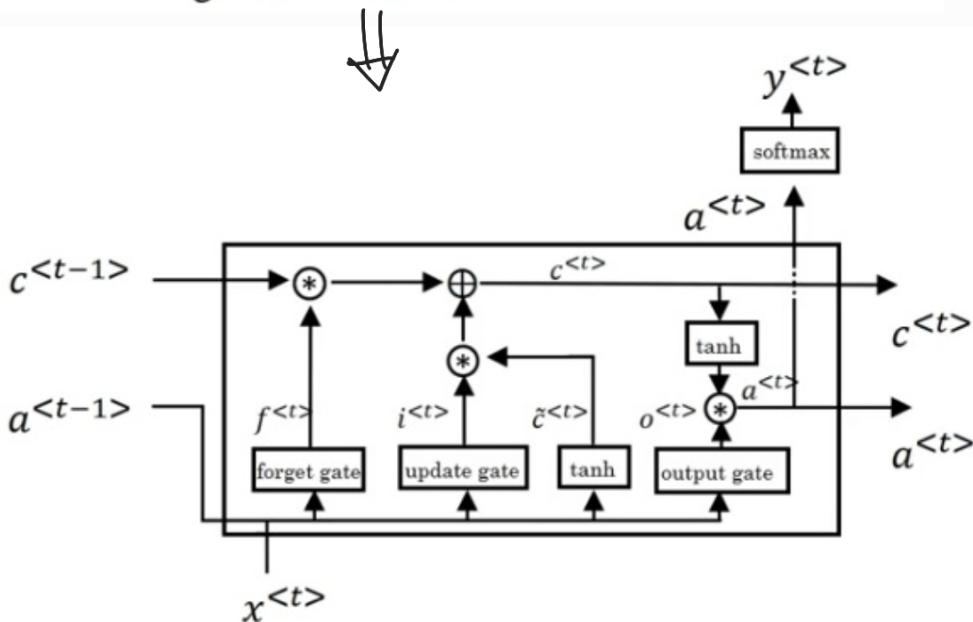
Update: $\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$

Forget: $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

Output: $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$



Variants:

- Peephole connection: for all gates, instead of $x^{<t>}$ use $c^{<t-1>}$, one-to-one.

What you should remember:

- An LSTM is similar to an RNN in that they both use hidden states to pass along information, but an LSTM also uses a cell state, which is like a long-term memory, to help deal with the issue of vanishing gradients
- An LSTM cell consists of a cell state, or long-term memory, a hidden state, or short-term memory, along with 3 gates that constantly update the relevancy of its inputs:
 - A **forget** gate, which decides which input units should be remembered and passed along. It's a tensor with values between 0 and 1.
 - If a unit has a value close to 0, the LSTM will "forget" the stored state in the previous cell state.
 - If it has a value close to 1, the LSTM will mostly remember the corresponding value.
 - An **update** gate, again a tensor containing values between 0 and 1. It decides on what information to throw away, and what new information to add.
 - When a unit in the update gate is close to 1, the value of its candidate is passed on to the hidden state.
 - When a unit in the update gate is close to 0, it's prevented from being passed onto the hidden state.
 - And an **output** gate, which decides what gets sent as the output of the time step

2 - Long Short-Term Memory (LSTM) Network

The following figure shows the operations of an LSTM cell:

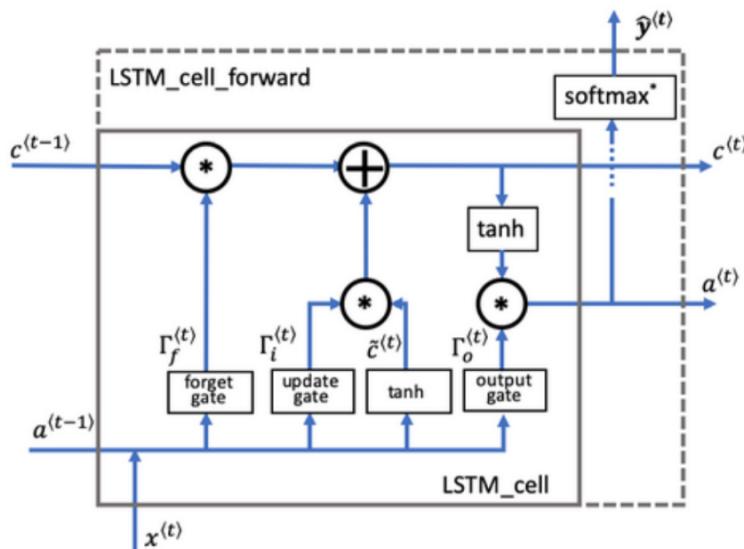


Figure 4: LSTM cell. This tracks and updates a "cell state," or memory variable $c^{(t)}$ at every time step, which can be different from $a^{(t)}$. Note, the `softmax` includes a dense layer and softmax.

Similar to the RNN example above, you'll begin by implementing the LSTM cell for a single time step. Then, you'll iteratively call it from inside a "for loop" to have it process an input with T_x time steps.

Overview of gates and states

Forget gate Γ_f

- Let's assume you are reading words in a piece of text, and plan to use an LSTM to keep track of grammatical structures, such as whether the subject is singular ("puppy") or plural ("puppies").
- If the subject changes its state (from a singular word to a plural word), the memory of the previous state becomes outdated, so you'll "forget" that outdated state.
- The "forget gate" is a tensor containing values between 0 and 1.
 - If a unit in the forget gate has a value close to 0, the LSTM will "forget" the stored state in the corresponding unit of the previous cell state.
 - If a unit in the forget gate has a value close to 1, the LSTM will mostly remember the corresponding value in the stored state.

Equation

$$\Gamma_f^{(t)} = \sigma(\mathbf{W}_f[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_f)$$

Explanation of the equation:

- \mathbf{W}_f contains weights that govern the forget gate's behavior.
- The previous time step's hidden state $[\mathbf{a}^{(t-1)}$ and current time step's input $\mathbf{x}^{(t)}$] are concatenated together and multiplied by \mathbf{W}_f .
- A sigmoid function is used to make each of the gate tensor's values $\Gamma_f^{(t)}$ range from 0 to 1.
- The forget gate $\Gamma_f^{(t)}$ has the same dimensions as the previous cell state $c^{(t-1)}$.
- This means that the two can be multiplied together, element-wise.
- Multiplying the tensors $\Gamma_f^{(t)} * c^{(t-1)}$ is like applying a mask over the previous cell state.
- If a single value in $\Gamma_f^{(t)}$ is 0 or close to 0, then the product is close to 0.
 - This keeps the information stored in the corresponding unit in $c^{(t-1)}$ from being remembered for the next time step.
- Similarly, if one value is close to 1, the product is close to the original value in the previous cell state.
 - The LSTM will keep the information from the corresponding unit of $c^{(t-1)}$, to be used in the next time step.

Variable names in the code

The variable names in the code are similar to the equations, with slight differences.

- `Wf` : forget gate weight \mathbf{W}_f
- `bf` : forget gate bias \mathbf{b}_f
- `ft` : forget gate $\Gamma_f^{(t)}$

Candidate value $\tilde{c}^{(t)}$

- The candidate value is a tensor containing information from the current time step that **may** be stored in the current cell state $c^{(t)}$.
- The parts of the candidate value that get passed on depend on the update gate.
- The candidate value is a tensor containing values that range from -1 to 1.
- The tilde "~" is used to differentiate the candidate $\tilde{c}^{(t)}$ from the cell state $c^{(t)}$.

Equation

$$\tilde{c}^{(t)} = \tanh(\mathbf{W}_c[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_c)$$

Explanation of the equation

- The \tanh function produces values between -1 and 1.

Variable names in the code

- `cct` : candidate value $\tilde{c}^{(t)}$

Update gate Γ_i

- You use the update gate to decide what aspects of the candidate $\tilde{c}^{(t)}$ to add to the cell state $c^{(t)}$.
- The update gate decides what parts of a "candidate" tensor $\tilde{c}^{(t)}$ are passed onto the cell state $c^{(t)}$.
- The update gate is a tensor containing values between 0 and 1.
 - When a unit in the update gate is close to 1, it allows the value of the candidate $\tilde{c}^{(t)}$ to be passed onto the hidden state $c^{(t)}$
 - When a unit in the update gate is close to 0, it prevents the corresponding value in the candidate from being passed onto the hidden state.
- Notice that the subscript "i" is used and not "u", to follow the convention used in the literature.

Equation

$$\Gamma_i^{(t)} = \sigma(\mathbf{W}_i[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_i)$$

Explanation of the equation

- Similar to the forget gate, here $\Gamma_i^{(t)}$, the sigmoid produces values between 0 and 1.
- The update gate is multiplied element-wise with the candidate, and this product ($\Gamma_i^{(t)} * \tilde{c}^{(t)}$) is used in determining the cell state $c^{(t)}$.

Variable names in code (Please note that they're different than the equations)

In the code, you'll use the variable names found in the academic literature. These variables don't use "u" to denote "update".

- `Wi` is the update gate weight \mathbf{W}_i (not "Wu")
- `bi` is the update gate bias \mathbf{b}_i (not "bu")
- `it` is the update gate $\Gamma_i^{(t)}$ (not "ut")

Cell state $c^{(t)}$

- The cell state is the "memory" that gets passed onto future time steps.
- The new cell state $c^{(t)}$ is a combination of the previous cell state and the candidate value.

Equation

$$c^{(t)} = \Gamma_f^{(t)} * c^{(t-1)} + \Gamma_i^{(t)} * \tilde{c}^{(t)}$$

Explanation of equation

- The previous cell state $c^{(t-1)}$ is adjusted (weighted) by the forget gate $\Gamma_f^{(t)}$
- and the candidate value $\tilde{c}^{(t)}$, adjusted (weighted) by the update gate $\Gamma_i^{(t)}$

Variable names and shapes in the code

- `c` : cell state, including all time steps, c shape (n_a, m, T_x)
- `c_next` : new (next) cell state, $c^{(t)}$ shape (n_a, m)
- `c_prev` : previous cell state, $c^{(t-1)}$, shape (n_a, m)

Output gate Γ_o

- The output gate decides what gets sent as the prediction (output) of the time step.
- The output gate is like the other gates, in that it contains values that range from 0 to 1.

Equation

$$\Gamma_o^{(t)} = \sigma(\mathbf{W}_o[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_o)$$

Explanation of the equation

- The output gate is determined by the previous hidden state $\mathbf{a}^{(t-1)}$ and the current input $\mathbf{x}^{(t)}$.
- The sigmoid makes the gate range from 0 to 1.

Variable names in the code

- \mathbf{W}_o : output gate weight, \mathbf{W}_o
- \mathbf{b}_o : output gate bias, \mathbf{b}_o
- \mathbf{o}_t : output gate, $\Gamma_o^{(t)}$

Hidden state $\mathbf{a}^{(t)}$

- The hidden state gets passed to the LSTM cell's next time step.
- It is used to determine the three gates ($\Gamma_f, \Gamma_u, \Gamma_o$) of the next time step.
- The hidden state is also used for the prediction $\mathbf{y}^{(t)}$.

Equation

$$\mathbf{a}^{(t)} = \Gamma_o^{(t)} * \tanh(\mathbf{c}^{(t)})$$

Explanation of equation

- The hidden state $\mathbf{a}^{(t)}$ is determined by the cell state $\mathbf{c}^{(t)}$ in combination with the output gate Γ_o .
- The cell state state is passed through the `tanh` function to rescale values between -1 and 1.
- The output gate acts like a "mask" that either preserves the values of $\tanh(\mathbf{c}^{(t)})$ or keeps those values from being included in the hidden state $\mathbf{a}^{(t)}$.

Variable names and shapes in the code

- \mathbf{a} : hidden state, including time steps. \mathbf{a} has shape (n_a, m, T_x)
- \mathbf{a}_{prev} : hidden state from previous time step. $\mathbf{a}^{(t-1)}$ has shape (n_a, m)
- \mathbf{a}_{next} : hidden state for next time step. $\mathbf{a}^{(t)}$ has shape (n_a, m)

Prediction $\mathbf{y}_{\text{pred}}^{(t)}$

- The prediction in this use case is a classification, so you'll use a softmax.

The equation is:

$$\mathbf{y}_{\text{pred}}^{(t)} = \text{softmax}(\mathbf{W}_y \mathbf{a}^{(t)} + \mathbf{b}_y)$$

Variable names and shapes in the code

- \mathbf{y}_{pred} : prediction, including all time steps. \mathbf{y}_{pred} has shape (n_y, m, T_x) . Note that $(T_y = T_x)$ for this example.
- $\mathbf{y}_{\text{t_pred}}$: prediction for the current time step t . $\mathbf{y}_{\text{t_pred}}^{(t)}$ has shape (n_y, m)

2.2 - Forward Pass for LSTM

Now that you have implemented one step of an LSTM, you can iterate this over it using a for loop to process a sequence of T_x inputs.

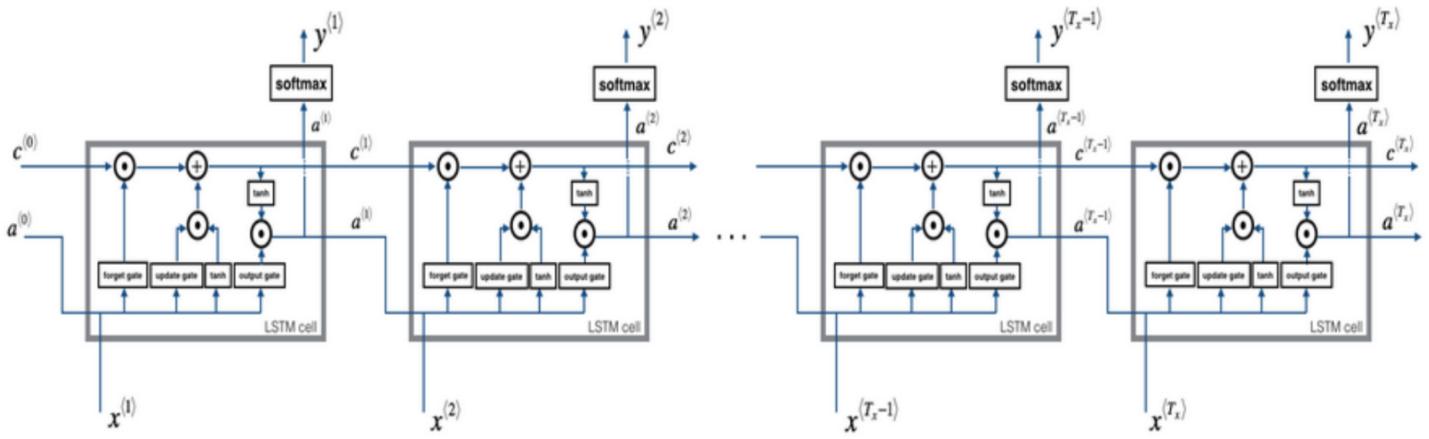


Figure 5: LSTM over multiple time steps.

Exercise 4 - lstm_forward

Implement `lstm_forward()` to run an LSTM over T_x time steps.

Instructions

- Get the dimensions n_x, n_a, n_y, m, T_x from the shape of the variables: `x` and `parameters`
- Initialize the 3D tensors a , c and y
 - a : hidden state, shape (n_a, m, T_x)
 - c : cell state, shape (n_a, m, T_x)
 - y : prediction, shape (n_y, m, T_x) (Note that $T_y = T_x$ in this example)
 - **Note** Setting one variable equal to the other is a "copy by reference". In other words, don't do ` $c = a$ `, otherwise both these variables point to the same underlying variable.
- Initialize the 2D tensor $a^{(t)}$
 - $a^{(t)}$ stores the hidden state for time step t . The variable name is `a_next`.
 - $a^{(0)}$, the initial hidden state at time step 0, is passed in when calling the function. The variable name is `a0`.
 - $a^{(t)}$ and $a^{(0)}$ represent a single time step, so they both have the shape (n_a, m)
 - Initialize $a^{(t)}$ by setting it to the initial hidden state ($a^{(0)}$) that is passed into the function.
- Initialize $c^{(t)}$ with zeros.
 - The variable name is `c_next`
 - $c^{(t)}$ represents a single time step, so its shape is (n_a, m)
 - **Note:** create `c_next` as its own variable with its own location in memory. Do not initialize it as a slice of the 3D tensor c . In other words, **don't** do