

Deep Learning

Chapter 8: Modern Recurrent Neural Networks

Mina Rezaei

Department of Statistics – LMU Munich

Winter Semester 2020



LECTURE OUTLINE

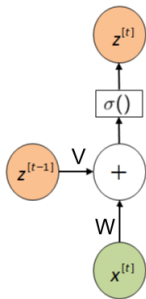
Long Short-Term Memory (LSTM)

Gated Recurrent Units (GRU)

Long Short-Term Memory (LSTM)

LONG SHORT-TERM MEMORY (LSTM)

The LSTM provides a way of dealing with vanishing gradients and modelling long-term dependencies.



- Until now, we simply computed

$$\mathbf{z}^{[t]} = \sigma(\mathbf{b} + \mathbf{V}^\top \mathbf{z}^{[t-1]} + \mathbf{W}^\top \mathbf{x}^{[t]})$$

LONG SHORT-TERM MEMORY (LSTM)

The LSTM provides a way of dealing with vanishing gradients and modelling long-term dependencies.

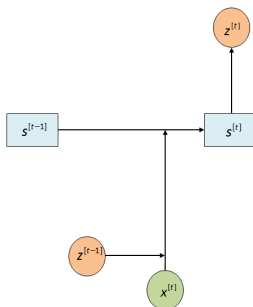


- Until now, we simply computed

$$\mathbf{z}^{[t]} = \sigma(\mathbf{b} + \mathbf{V}^\top \mathbf{z}^{[t-1]} + \mathbf{W}^\top \mathbf{x}^{[t]})$$

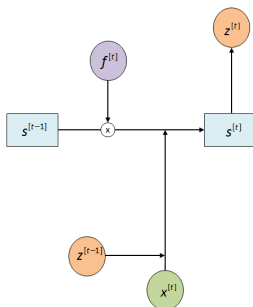
- Now we introduce the LSTM cell (which is a small network on its own).

LONG SHORT-TERM MEMORY (LSTM)



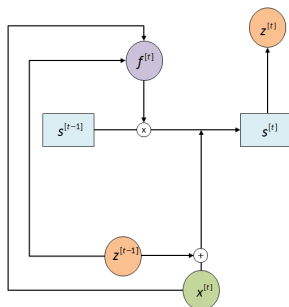
- The key to LSTMs is the **cell state** $s^{[t]}$.
- $s^{[t]}$ can be manipulated by different **gates** to forget old information, add new information, and read information out of it.
- Each gate is a vector of the same size as the $s^{[t]}$ with elements between 0 ("let nothing pass") and 1 ("let everything pass").

LONG SHORT-TERM MEMORY (LSTM)



- **Forget gate $f^{[t]}$** : indicates which information of the old cell state we should forget.
- Intuition: Think of a model trying to predict the next word based on all the previous ones. The cell state might include the gender of the present subject, so that the correct pronouns can be used. When we now see a new subject, we want to forget the gender of the old one.

LONG SHORT-TERM MEMORY (LSTM)

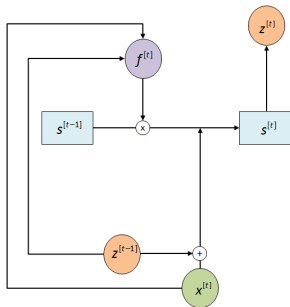


- We obtain the forget gate by computing

$$\mathbf{f}^{[t]} = \sigma(\mathbf{b}_f + \mathbf{V}_f^\top \mathbf{z}^{[t-1]} + \mathbf{W}_f^\top \mathbf{x}^{[t]})$$

- $\sigma()$ is a sigmoid, squashing the values to $[0, 1]$, and \mathbf{V}_f , \mathbf{W}_f are forget gate specific weights.

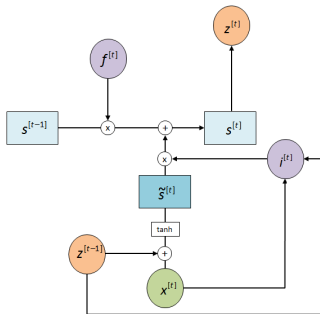
LONG SHORT-TERM MEMORY (LSTM)



- To compute the cell state $\mathbf{s}^{[t]}$, the first step is to multiply (element-wise) the previous cell state $\mathbf{s}^{[t-1]}$ by the forget gate $\mathbf{f}^{[t]}$.

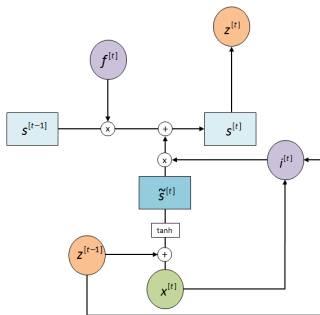
$$\mathbf{f}^{[t]} \odot \mathbf{s}^{[t-1]}, \text{ with } \mathbf{f}^{[t]} \in [0, 1]$$

LONG SHORT-TERM MEMORY (LSTM)



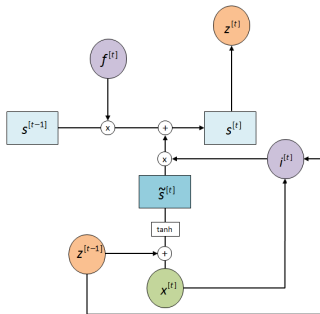
- **Input gate** $i^{[t]}$: indicates which new information should be added to $s^{[t]}$.
- Intuition: In our example, this is where we add the new information about the gender of the new subject.

LONG SHORT-TERM MEMORY (LSTM)



- The new information is given by
$$\tilde{s}^{[t]} = \tanh(\mathbf{b} + \mathbf{V}^\top \mathbf{z}^{[t-1]} + \mathbf{W}^\top \mathbf{x}^{[t]}) \in [-1, 1].$$
- The input gate is given by
$$i^{[t]} = \sigma(\mathbf{b}_i + \mathbf{V}_i^\top \mathbf{z}^{[t-1]} + \mathbf{W}_i^\top \mathbf{x}^{[t]}) \in [0, 1].$$
- \mathbf{W} and \mathbf{V} are weights of the new information, \mathbf{W}_i and \mathbf{V}_i the weights of the input gate.

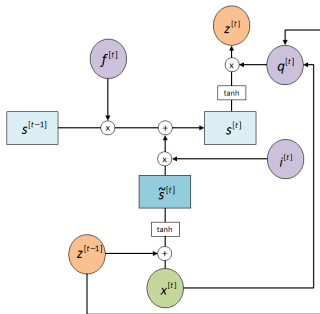
LONG SHORT-TERM MEMORY (LSTM)



- Now we can finally compute the cell state $\mathbf{s}^{[t]}$:

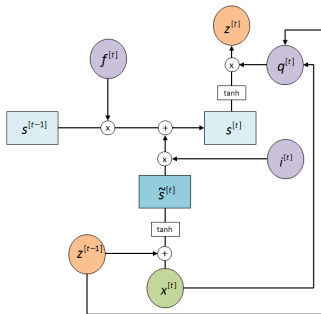
$$\mathbf{s}^{[t]} = \mathbf{f}^{[t]} \odot \mathbf{s}^{[t-1]} + \mathbf{i}^{[t]} \odot \tilde{\mathbf{s}}^{[t]}$$

LONG SHORT-TERM MEMORY (LSTM)



- **Output gate $q^{[t]}$:** Indicates which information from the cell state is filtered.
- It is given by $q^{[t]} = \sigma(\mathbf{b}_q + \mathbf{V}_q^\top \mathbf{z}^{[t-1]} + \mathbf{W}_q^\top \mathbf{y}^{[t]})$, with specific weights $\mathbf{W}_q, \mathbf{V}_q$.

LONG SHORT-TERM MEMORY (LSTM)



- Finally, the new state $z^{[t]}$ of the LSTM is a function of the cell state, multiplied by the output gate:

$$z^{[t]} = q^{[t]} \odot \tanh(s^{[t]})$$

Gated Recurrent Units (GRU)

GATED RECURRENT UNITS (GRU)

- The key distinction between regular RNNs and GRUs is that the latter support gating of the hidden state.
- Here, we have dedicated mechanisms for when a hidden state should be updated and also when it should be reset.
- These mechanisms are learned and they address the listed concerns:
 - avoid the vanishing/exploding gradient problem which comes with a standard recurrent neural network.
 - GRUs are able to solve the vanishing gradient problem by using an update gate and a reset gate.
 - the update gate controls information that flows into memory, and the reset gate controls the information that flows out of memory.

RESET GATES AND UPDATE GATES

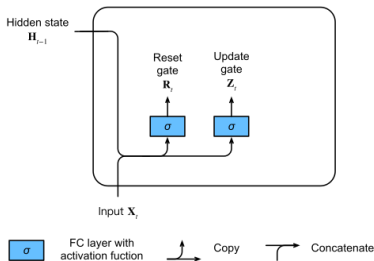


Figure: Reset and update gate in a GRU.

- as shown by Figure1, given the current time step input X_t and the hidden state of the previous time step H_{t-1} . The output is given by a fully connected layer with a sigmoid as its activation function.

RESET GATES AND UPDATE GATES

- For a given time step t , the minibatch input is $\mathbf{X}_t \in \mathbb{R}^{n \times d}$ (number of examples: n , number of inputs: d) and the hidden state of the last time step is $\mathbf{H}_{t-1} \in \mathbb{R}^{n \times h}$ (number of hidden states: h). Then, the reset gate $\mathbf{R}_t \in \mathbb{R}^{n \times h}$ and update gate $\mathbf{Z}_t \in \mathbb{R}^{n \times h}$ are computed as follows:
- $\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r)$
- $\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z)$
- Here, $\mathbf{W}_{xr}, \mathbf{W}_{xz} \in \mathbb{R}^{d \times h}$ and $\mathbf{W}_{hr}, \mathbf{W}_{hz} \in \mathbb{R}^{h \times h}$ are weight parameters and $\mathbf{b}_r, \mathbf{b}_z \in \mathbb{R}^{1 \times h}$ are biases. We use a sigmoid function to transform input values to the interval $(0, 1)$.

RESET GATES IN ACTION

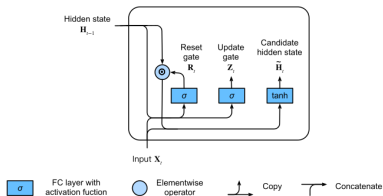


Figure: Candidate hidden state computation in a GRU. The multiplication is carried out elementwise.

- We begin by integrating the reset gate with a regular latent state updating mechanism. In a conventional RNN, we would have an hidden state update of the form:
- $\mathbf{H}_t = \tanh(\mathbf{x}_t \mathbf{W}_{xh} + \mathbf{H}_{t-1} \mathbf{W}_{hh} + \mathbf{b}_h)$.
- $\tilde{\mathbf{H}}_t = \tanh(\mathbf{x}_t \mathbf{W}_{xh} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{hh} + \mathbf{b}_h)$.

RESET GATES IN ACTION

- Figure 2 illustrates the computational flow after applying the reset gate.

UPDATE GATES IN ACTION

- aa
- Next we need to incorporate the effect of the update gate Z_t .
- The gating variable Z_t can be used for this purpose, simply by taking elementwise convex combinations between both candidates.
- This leads to the final update equation for the GRU.
- $\mathbf{H}_t = \mathbf{Z}_t \odot \mathbf{H}_{t-1} + (1 - \mathbf{Z}_t) \odot \tilde{\mathbf{H}}_t$.

UPDATE GATES IN ACTION

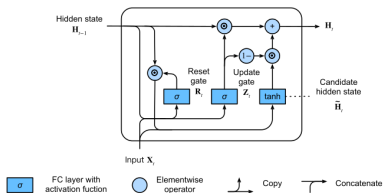


Figure: Hidden state computation in a GRU. As before, the multiplication is carried out elementwise.

These designs can help us cope with the vanishing gradient problem in RNNs and better capture dependencies for time series with large time step distances. In summary, GRUs have the following two distinguishing features:

- Reset gates help capture short-term dependencies in time series.
- Update gates help capture long-term dependencies in time series.

REFERENCES



Ian Goodfellow, Yoshua Bengio and Aaron Courville (2016)

Deep Learning

<http://www.deeplearningbook.org/>



Michael Nguyen (2018)

Illustrated Guide to LSTM's and GRU's: A step by step explanation newblock

<https://towardsdatascience.com/>

[illustrated-guide-to-lstms-and-gru-s-a-step-by-step-explanation-44e](#)