Ali Athar <athar@vision.rwth-aachen.de>
Sabarinath Mahadevan <mahadevan@vision.rwth-aachen.de>

# Exercise 6: Recurrent Neural Networks

due **before** 2020-01-29

**Important information regarding the exercises:**

- The exercises are not mandatory. Still, we strongly encourage you to solve them! All submissions will be corrected. If you submit your solution, please read on:

- Send your submission via mail with the subject `[ML] exercise solution 6` to <**athar@vision.rwth-aachen.de**> and <**mahadevan@vision.rwth-aachen.de**>.

- Due to the large number of participants, we require you to submit your solution to L$^2$P **in groups of 3 to 4 students**. You can use the **Discussion Forum** on L$^2$P to organize groups.

- If applicable submit your code solution as a zip/tar.gz file named `mn1_mn2_mn3.{zip/tar.gz}` with your **matriculation numbers** (`mn`).

- Please do **not** include the data files in your submission!

- Please upload your pen & paper problems as PDF. Alternatively, you can also take pictures (.png or .jpeg) of your hand written solutions. Please make sure your handwriting is legible, the pictures are not blurred and taken under appropriate lighting conditions. All non-readable submissions will be discarded immediately.

**Question 1: Basic RNN model and BPTT** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . $(\mathbf{\Sigma = 9})$

This exercise deals with the implementation of a basic RNN and LSTM cell. To this end, you are given two code frameworks which will be completed in the following tasks. The first framework will introduce you to the challenges of implementing a basic RNN from scratch by providing guidance for the key steps in the process, namely forward propagation and backward propagation through time (BPTT). This will be followed by a test case, where the network is trained to remember an arbitrary given sequence across several time steps, and outputting the same sequence after reading a specific symbol. This kind of application tests the memory capabilities of the RNN, as well as the ability to learn the concept of first remembering and then repeating.

The second framework is implemented in TensorFlow. This implementation is based on the abstract class `RNNCell` that defines the basic interface for our RNNs. Your task will be, to complete the interface implementation for the child class `MyLSTMCell`.

Specifically, this comprises implementing a `__call__` method for the LSTM cell which basically adds an unrolled RNN to TensorFlow's graph that runs on a given input. The outputs are fed into a linear layer on top of the cell in order to map the cell output to the desired dimension. The parameters of this layer are shared among all steps. Finally, the LSTM cell's performance will be evaluated on two test cases.

The first one will train the network to continue a sine curve of an unknown frequency after being presented a few initial samples. To this end, ground truth samples from training sine curves are presented to the cell which in every time step should then predict the next sample using the linear layer on top. The prediction error is measured by a mean squared error over the entire training sequence.

The second one is similar to the test case for the from scratch implementation. The only difference is that the sequence is separated from the delimiter by a certain number of spaces. By doing this the complexity of the task is slightly increased and this will demonstrate the superiority of LSTMs in memory retention tasks, with regards to basic RNNs. After training, the network will be tested on previously unseen data and evaluated according to an accuracy measure.

**Requirements:** The following packages are required: `numpy`, `tensorflow` and `matplotlib`. Please install any missing packages using `pip install "package"` (without colons) from the command line.

(a) **Forward propagation of basic RNN:** As a first task, you are asked to complete **(2 pts)** the code of the `BasicRNNCell.fprop` function. Remember that the equations for $h_t$, the hidden state at time $t$ are given by:

$$\mathbf{h}_{t+1} = \text{ReLU}\left(\mathbf{W}_{hh} \cdot \mathbf{h}_t + \mathbf{W}_{hx} \cdot \mathbf{x}_t + \mathbf{b}_h\right)$$

Where $\mathbf{W}$ are the weight matrices and $\mathbf{x}$ the input. Take note, that for simplicities sake, $\mathbf{h}_0$ is always initialized to zero. In a real application environment, the state of $\mathbf{h}_0$ would be an additional parameter that can be passed to the function. The linear mapping to the output variable $y$

$$\mathbf{y}_t = \mathbf{W}_{yh} \cdot \mathbf{h}_t$$

is done using the separate class `LinearLayer` so that multiple layers of `BasicRNNCells` may be stacked on top of each other.

*Hint:* For information about the interface you may want to take a look at the implementation for the linear layer already finished in `LinearLayer.py` as well as `BasicRNNCell.__init__` for the variable names.

(b) **BPTT of basic RNN:** Now you are asked to implement BPTT by completing the **(3 pts)** code in `BasicRNNCell.bprop`. In the following, $\hat{\mathbf{h}}_t$ will denote the derivative with respect to the hidden state at time $t$ before the ReLU is applied. The BPTT steps for the forward-propagation above can be written as follows:

$$\partial \mathbf{h}_t = W_{yh}^T \partial \mathbf{y}_t + \mathbf{W}_{hh}^T \cdot \partial \hat{\mathbf{h}}_{t+1}$$
$$\partial \hat{\mathbf{h}}_t = \text{ReLU}_t\left(\partial \mathbf{h}_t\right)$$
$$\partial \mathbf{W}_{hh} += \partial \hat{\mathbf{h}}_t \cdot \mathbf{h}_{t-1}^T$$
$$\partial \mathbf{W}_{hx} += \partial \hat{\mathbf{h}}_t \cdot \mathbf{x}_t^T$$
$$\partial \mathbf{b}_h += \partial \hat{\mathbf{h}}_t$$
$$\partial \mathbf{x}_t = \mathbf{W}_{hx}^T \cdot \partial \hat{\mathbf{h}}_t$$

The $+=$ is necessary in order to accumulate the gradient over time. It is important not to calculate the ReLU function on the derivatives, but to apply the same ReLU function, that was previously applied in timestep $t$. This means, that you need to remember which values were set to 0 in all timesteps $t$. In the backpropagation you need to set the same values of $\partial \hat{\mathbf{h}}_t = 0$ as these values cannot pass down a derivative, since they were set to zero in the forward-propagation. It is also of note, that in the code the derivatives dys which are given as an input to the bprop function are the derivatives that are passed down from higher layers, i.e. the linear layer, and have already been multiplied with $W_{yh}^T$, i.e. $dys_t$ will be the result of $W_{yh}^T \cdot \partial \mathbf{y}_t$. Additionally,

$\partial\hat{\mathbf{h}}_{t+1}$ should be initialized to zero for $t = t_{\text{end}}$. Once you have completed implementing the forward- and backward-propagation, you can use the provided `run_check_grads` function to verify, that both are working as intended.

*Hint:* For information about the interface you may want to take a look at the implementation for the linear layer already finished in `LinearLayer.py`.

(c) **Exploding gradient:** As was already introduced in the lecture, one big flaw of **(0 pts)** recurrent networks, is that the gradient of parameter matrices may vanish or become extremely large. After you have completed the previous tasks, use the function `run_exploding_grads_test` to visualize how quickly the mean absolute value of the gradient can increase. You may also try out different numbers for the hidden layer, sequence length and input vocabulary to see how this affects the magnitude of the gradient.

(d) **Gradient clipping:** As a method for mitigation of exploding gradients, gradient **(1 pt)** clipping was introduced. This method computes the norm of the gradient after BPTT and scales it to a fixed value. The pseudo-code for the algorithm is presented in the following

```
1    if ||grad|| > threshhold then
2        grad ← threshhold/||grad|| grad
```

Your task is to implement `common.clip_gradients`. After you finished your implementation use `run_exploding_grads_test` to make sure that the gradient clipping has taken effect correctly.

*Hint:* In order to make sure that the changes take effect in the gradient, which is passed as an argument to the function, you should use the `/=` or `*=` assignment operators.

(e) **Memorization Task:** Finally, we want our RNN to learn to perform a simple memo- **(2 pts)** rization task. The goal is to take an input sequence and read until a specific character is encountered. After this point the RNN should output all the characters that were previously read in the correct order. The input sequence will look something like

$$\begin{pmatrix} c_1 & \cdots & c_{\text{to\_remember\_len}} & \# & \underbrace{\sqcup \quad \cdots \quad \sqcup}_{\text{to\_remember\_len } times} \end{pmatrix}$$

whereas the target output sequence will look like

$$\begin{pmatrix} \underbrace{\sqcup \quad \cdots \quad \sqcup}_{\text{to\_remember\_len } times} & \# & c_1 & \cdots & c_{\text{to\_remember\_len}} \end{pmatrix}$$

To this end you are required to implement the function `run_memory_test` which computes the accuracy of the network on a randomly generated test set of data. For data generation you may want to use the function `generate_data_memory`, which returns data and label matrices. Additionally, the function `run_forward_pass` may be helpful to you in order to obtain the soflty assigned classes. Remember, that the only interesting values for the accuracy measure are the ones generated after the delimiter has been read. Afterwards you may try to adjust the command line parameters (e.g. learning rate, number of steps) so that the test accuracy is maximized.

(f) **Vanishing gradient:** One method to fix the vanishing gradient issue that was pre- **(1 pt)** sented, is to initialize $\mathbf{W}_{hh}$ to the identity matrix and to use ReLU non-linearities. As the latter is already implemented in the code, your task is to change the initialization of $\mathbf{W}_{hh}$ to the identity matrix in `BasicRNNCell.__init__`. Do you notice any change in the accuracy of the memorization task?

**Question 2: LSTM in Tensorflow** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . $\left(\Sigma = \mathbf{10}\right)$

(a) **Implementing LSTM:** Now, we will implement our version of an LSTM cell as presented in the lecture. Start by adding the variables to `MyLSTMCell.__init__`. **(2 pts)**
*Hint:* You may want to take a look at the implementation in `BasicRNNCell.__init__` for a point of reference.

(b) As for the basic RNN cell, implement the `MyLSTMCell.__call__` function. **(5 pts)**
The input **x** is a tensor with dimensions as follows:
- 0: Batch
- 1: Time steps
- 2: Input value for the batch and time step.

The function should return a list of output tensors corresponding to the cell output at every time step and the hidden state tensor $\mathbf{h}_T$ after the last time step T.
*Hint:* The state is described by a `MyLSTMState` tuple containing hidden and cell state, respectively (see `MyLSTMCell.py`, line 5). You may also look at the implementation of `BasicRNNCell` for reference.

(c) As mentioned above, we want our cell to generate a sine curve within a range of frequencies. Therefore, implement the mean squared error calculation in `train.` **(1 pt)**
`get_loss_sin`.
*Hint:* The ground truth is a 2D tensor where rows are used for the batch and columns correspond to sine evaluation steps.
*Hint:* If `RNNCell.__call__` is implemented correctly, the $i$th column is the ground truth value for the cell prediction at the $i$th step. No shifting needed.

(d) Finally, implement `train.get_train_op` which should return a gradient descent **(1 pt)**
based training operation. Remember that in general gradient clipping may be useful to cope with exploding gradients. You can now test the sine curve generation by setting the command line parameter --task 2. In this setup, the first 30 samples of a sin curve generated with a random frequency and phase shift are given as an input to the network and afterwards the network extrapolates the curve by feeding back its own predictions as its input. The results are then written out as png files in the folder of the source code. One possible configuration using the basic RNN is set as default. To use the LSTM, set the command line parameter --RNN LSTM. If you try the LSTM, you may want to have a fairer comparison in terms of number of parameters. How can you achieve this?
*Hint:* You may want to use `tf.clip_by_global_norm` and a norm of 5.0

(e) Now, we want to challenge the memory capacities of our basic RNN cell and our **(1 pt)**
implementation of the LSTM. To this end, you should implement a memory task which tests the length of past to present dependencies that the cell can learn. Therefore, the cell is trained to run on the sequence

$$\left(r_1 \quad \cdots \quad r_{\text{to\_remember\_len}} \quad \underbrace{\text{␣} \quad \cdots \quad \text{␣}}_{\text{blank\_separation\_len } times} \quad \# \quad r_1 \quad \cdots \quad r_{\text{to\_remember\_len}}\right)$$

with $r_1, \ldots, r_{\text{to\_remember\_len}}$ chosen from $\{0, \ldots, 8\}$ uniformly at random while generating the output

$$\left( \quad \underbrace{\text{␣} \quad \cdots \quad \text{␣}}_{\text{to\_remember\_len+blank\_separation\_len+1 } times} \quad r_1 \quad \cdots \quad r_{\text{to\_remember\_len}}\right) \cdot$$

Thus, the challenge is to remember the sequence at the beginning over a certain number of blanks and to start recalling after the delimiter.

First, complete `train.get_loss_memory` using an appropriate loss and return a single scalar loss value obtained by summing over the individual losses for each element of the batch.

(f) Implement the method `test.run_memory_test` which evaluates the memory capability of the trained cell. After running the test iterator initializer you can evaluate the given loss on the test dataset. Besides, it may be more interesting to calculate the accuracy of recalling the sequence which should be remembered. Therefore, also calculate the accuracy for the last to_remember_len outputs. Which model does perform better?

*Hint:* See `train.run_training` to see how the test data can be consumed.

*Hint:* You might want to try to increase the number of training epochs for this task.