

# Homework 2: Convolutional Neural Networks and Recurrent Neural Networks

DS-GA 1008 Deep Learning

Spring 2024

The goal of homework 2 is to get you to work with convolutional neural networks and recurrent neural networks.

In the theoretical part, you will work on figuring out how backpropagation works in these networks. In part 2, we will implement and train them.

In part 1, you should submit all your answers in a pdf file. As before, we recommend using  $\text{\LaTeX}$ .

For part 2, you will implement some neural networks by adding your code to the provided ipynb file.

As before, please use numerator layout.

Submit the following files in a zip file `your_net_id.zip` through NYU Brightspace:

- `hw2_theory.pdf`
- `hw2_moe.ipynb`
- `hw2_cnn.ipynb`
- `hw2_rnn.ipynb`
- `08-seq_classification.ipynb`

The following behaviors will result in penalty of your final score:

1. 10% penalty for submitting your file without using the correct naming format (including naming the zip file, PDF file or python file wrong, adding extra files in the zip folder, like the testing scripts in your zip file).
2. 20% penalty for code submission that cannot be executed following the steps we mentioned.

# 1 Theory (50pt)

## 1.1 Convolutional Neural Networks (15 pts)

- (a) (1 pts) Given an input image of dimension  $11 \times 21$ , what will be output dimension after applying a convolution with  $5 \times 4$  kernel, stride of 4, and no padding?
- (b) (2 pts) Given an input of dimension  $C \times H \times W$ , what will be the dimension of the output of a convolutional layer with kernel of size  $K \times K$ , padding  $P$ , stride  $S$ , dilation  $D$ , and  $F$  filters. Assume that  $H \geq K$ ,  $W \geq K$ .
- (c) (12 pts) In this section, we are going to work with 1-dimensional convolutions. Discrete convolution of 1-dimensional input  $x[n]$  and kernel  $k[n]$  is defined as follows:

$$s[n] = (x * k)[n] = \sum_m x[n - m]k[m]$$

However, in machine learning convolution is usually implemented as cross-correlation, which is defined as follows:

$$s[n] = (x * k)[n] = \sum_m x[n + m]k[m]$$

Note the difference in signs, which will get the network to learn an “flipped” kernel. In general it doesn’t change much, but it’s important to keep it in mind. In convolutional neural networks, the kernel  $k[n]$  is usually 0 everywhere, except a few values near 0:  $\forall_{|n|>M} k[n] = 0$ . Then, the formula becomes:

$$s[n] = (x * k)[n] = \sum_{m=-M}^M x[n + m]k[m]$$

Let’s consider an input  $x[n] \in \mathbb{R}^5$ , with  $1 \leq n \leq 7$ , e.g. it is a length 7 sequence with 5 channels. We consider the convolutional layer  $f_W$  with one filter, with kernel size 3, stride of 2, no dilation, and no padding. The only parameters of the convolutional layer is the weight  $W$ ,  $W \in \mathbb{R}^{1 \times 5 \times 3}$ , there’s no bias and no non-linearity.

- (i) (1 pts) What is the dimension of the output  $f_W(x)$ ? Provide an expression for the value of elements of the convolutional layer output  $f_W(x)$ . Example answer format here and in the following sub-problems:  $f_W(x) \in \mathbb{R}^{42 \times 42 \times 42}$ ,  $f_W(x)[i, j, k] = 42$ .
- (ii) (3 pts) What is the dimension of  $\frac{\partial f_W(x)}{\partial W}$ ? Provide an expression for the values of the derivative  $\frac{\partial f_W(x)}{\partial W}$ .
- (iii) (3 pts) What is the dimension of  $\frac{\partial f_W(x)}{\partial x}$ ? Provide an expression for the values of the derivative  $\frac{\partial f_W(x)}{\partial x}$ .

- (iv) (5 pts) Now, suppose you are given the gradient of the loss  $\ell$  w.r.t. the output of the convolutional layer  $f_W(x)$ , i.e.  $\frac{\partial \ell}{\partial f_W(x)}$ . What is the dimension of  $\frac{\partial \ell}{\partial W}$ ? Provide an expression for  $\frac{\partial \ell}{\partial W}$ . Explain similarities and differences of this expression and expression in (i).

## 1.2 Recurrent Neural Networks (30 pts)

### 1.2.1 Part 1

In this section we consider a simple recurrent neural network defined as follows:

$$c[t] = \sigma(W_c x[t] + W_h h[t-1]) \quad (1)$$

$$h[t] = c[t] \odot h[t-1] + (1 - c[t]) \odot W_x x[t] \quad (2)$$

where  $\sigma$  is element-wise sigmoid,  $x[t] \in \mathbb{R}^n$ ,  $h[t] \in \mathbb{R}^m$ ,  $W_c \in \mathbb{R}^{m \times n}$ ,  $W_h \in \mathbb{R}^{m \times m}$ ,  $W_x \in \mathbb{R}^{m \times n}$ ,  $\odot$  is Hadamard product,  $h[0] \doteq 0$ .

- (a) (4 pts) Draw a diagram for this recurrent neural network, similar to the diagram of RNN we had in class. We suggest using [diagrams.net](https://diagrams.net).
- (b) (1pts) What is the dimension of  $c[t]$ ?
- (c) (5 pts) Suppose that we run the RNN to get a sequence of  $h[t]$  for  $t$  from 1 to  $K$ . Assuming we know the derivative  $\frac{\partial \ell}{\partial h[t]}$ , provide dimension of and an expression for values of  $\frac{\partial \ell}{\partial W_x}$ . What are the similarities of backward pass and forward pass in this RNN?
- (d) (2pts) Can this network be subject to vanishing or exploding gradients? Why?

### 1.2.2 Part 2

We define an AttentionRNN(2) as

$$q_0[t], q_1[t], q_2[t] = Q_0 x[t], Q_1 h[t-1], Q_2 h[t-2] \quad (3)$$

$$k_0[t], k_1[t], k_2[t] = K_0 x[t], K_1 h[t-1], K_2 h[t-2] \quad (4)$$

$$v_0[t], v_1[t], v_2[t] = V_0 x[t], V_1 h[t-1], V_2 h[t-2] \quad (5)$$

$$w_i[t] = q_i[t]^\top k_i[t] \quad (6)$$

$$a[t] = \text{softargmax}([w_0[t], w_1[t], w_2[t]]) \quad (7)$$

$$h[t] = \sum_{i=0}^2 a_i[t] v_i[t] \quad (8)$$

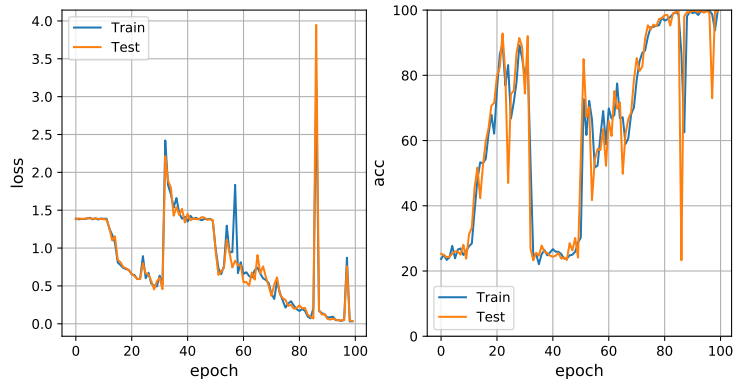
Where  $x[t], h[t] \in \mathbb{R}^n$ , and  $Q_i, K_i, V_i \in \mathbb{R}^{n \times n}$ . We define  $h[t] = \mathbf{0}$  for  $t < 1$ . You may safely ignore these bases cases in the following questions.

- (a) (4 pts) Draw a diagram for this recurrent neural network
- (b) (1 pt) What is the dimension of  $a[t]$ ?
- (c) (3 pts) Extend this to,  $\text{AttentionRNN}(k)$ , a network that uses the last  $k$  state vectors  $h$ . Write out the system of equations that defines it. You may use set notation or ellipses (...) in your definition.
- (d) (3 pts) Modify the above network to produce  $\text{AttentionRNN}(\infty)$ , a network that uses every past state vector. Write out the system of equations that defines it. You may use set notation or ellipses (...) in your definition. HINT: We can do this by tying together some set of parameters, e.g. weight sharing.
- (e) (5 pts) Suppose the loss  $\ell$  is computed. Please write down the expression for  $\frac{\partial h[t]}{\partial h[t-1]}$  for  $\text{AttentionRNN}(2)$ .
- (f) (2 pts) Suppose we know the derivative  $\frac{\partial h[t]}{\partial h[T]}$ , and  $\frac{\partial \ell}{\partial h[t]}$  for all  $t > T$ . Please write down the expression for  $\frac{\partial \ell}{\partial h[T]}$  for  $\text{AttentionRNN}(k)$ .

### 1.3 Debugging loss curves (5pts)

When working with notebook

[08-seq\\_classification](#), we saw RNN training curves. In Section 8 “Visualize LSTM”, we observed some “kinks” in the loss curve.



1. (1pts) What caused the spikes on the left?
2. (1pts) How can they be higher than the initial value of the loss?
3. (1pts) What are some ways to fix them?
4. (2pts) Explain why the loss and accuracy are at these values before training starts. You may need to check the task definition in the notebook.

## 2 Implementation (50pts + 7pts extra credit)

There are three notebooks in the practical part:

- (10 + 2pts) Mixture of Experts: [hw2\\_moe.ipynb](#)
- (20pts) Convolutional Neural Networks notebook: [hw2\\_cnn.ipynb](#)
- (15pts) Recurrent Neural Networks notebook: [hw2\\_rnn.ipynb](#)
- (5pts + 5pts extra credit) : This builds on Section 1.3 of the theoretical part.
  - (5pts) Change the model training procedure of Section 8 in [08-seq\\_classification](#) to make the training curves have no spikes. You should only change the training of the model, and not the model itself or the random seed.
  - (5pts extra credit) Visualize the gradients and weights throughout training before and after you fix the training procedure.

**Plase use your NYU Google Drive account to access the notebooks.** First two notebooks contain parts marked as TODO, where you should put your code. These notebooks are Google Colab notebooks, you should copy them to your drive, add your solutions, and then download and submit them to NYU Brightspace. The notebook from the class, if needed, can be uploaded to Colab as well.