# Final Exam Long Answer Questions

Introduction to Deep Learning Simon Fraser University

<u>Instructions</u>: Below are listed 6 questions. Exactly 2 of these questions will appear on your exam. You can upload two files, one for each answer. You should explain your answers, even if not explicitly asked to do so. For example, if a question asks: "what is the derivative of equation E", you should not only give the derivative, but also outline how you obtained it.

The exam is open-book, in the sense that you can consult the textbook and on-line resources like Wikipedia. The university policy on academic dishonesty and plagiarism (cheating) will be taken very seriously in this course. *Everything submitted should be your own writing or coding.* You must not let other students copy your work. If I determine that you have copied, you will receive 0 marks.

<u>Group Work:</u> Discussing the problems is okay, for example to understand the concepts involved. If you work in a group, put down the name of all members of your group. There should be no group submissions. *Each group member should write up their own solution* to show their own understanding.

### Question on Recurrent Neural Networks

We want to process two binary input sequences with 0-1 entries and determine if they are equal. For notation, let  $x_1 = x_1^{(1)}, x_1^{(2)}, ..., x_1^{(T)}$  be the first input sequence and  $x_2 = x_2^{(1)}, x_2^{(2)}, ..., x_2^{(T)}$  be the second. We use the RNN architecture shown in the Figure.



The corresponding update equations are as follows.

<b>h</b> <sup>(t)</sup> =g( <b>Wx</b> <sup>(t)</sup> + <b>b</b> )	
$y^{(t)}=g(v^{T}h^{(t)}+ry^{(t-1)}+c)$	for t>1
$y^{(t)}=g(\mathbf{v}^{T}\mathbf{h}^{(t)}+c_{0})$	for t=1

Where  $\mathbf{v}^{\mathsf{T}}$  is the transpose of vector  $\mathbf{v}$  and the activation function g is defined as follows.

<i>g(z)</i> =1	for <i>z</i> > 0
<i>g(z)</i> =0	for $z \leq 0$

Described in words, the parameters are as follows.

W	2x2 weight matrix
b	2-dimensional bias vector
v	2-dimensional weight vector
r	Scalar recurrent weight
С	Scalar bias for all time steps except the first
С0	Scalar bias for the first time step

I suggest the following strategy for solving this problem.

- At time step *t*, the neural network is fed two inputs  $x_1^{(t)}$  and  $x_2^{(t)}$ .
- Use the two hidden units  $h_1^{(t)}$  and  $h_2^{(t)}$  to determine if the current inputs match.
- Use the output unit  $y^{(t)}$  to compute whether all inputs have matched up to the current time.

Specify parameter values that correctly implement this function, like in the table shown. (You do not have to write your answer in the table). Justify why you think your parameter values are correct.

W	Your solution
b	Your solution
v	Your solution
r	Your solution
С	Your solution
С0	Your solution

## Question on Gated Recurrent Neural Networks

Suppose we want to build a gated RNN cell <u>that sums its inputs over time</u>. What should be the gating values be? To focus on the gating aspect, your design can change the activation function of the RNN cell (e.g., replace tanh by linear).

- 1. For the LSTM architecture as explained in Section 4.6. of the text, what should be the value of the input gate and the forget gate?
- For the GRU architecture, what should be the value of the reset gate and the update gate? The GRU architecture is described in the slides and in on-line sources like this one <u>https://towardsdatascience.com/illustrated-guide-to-lstms-and-gru-s-a-step-by-stepexplanation-44e9eb85bf21</u>.

#### Question on Variational Auto-Encoders

- 1. Compare the loss function for the Variational Auto-Encoder (Figure 7.8 in the text) to the loss function for an associative auto-encoder (Figure 7.1 in the text). Which parts are similar and which are different?
- 2. How does the VAE architecture allow it to generate new data points, especially compared to associative auto-encoder, which cannot generate new data points?
- 3. Let *d* be the latent embedding dimension. The VAE encoder outputs a mean vector  $\boldsymbol{\mu} = (m_1, m_2, ..., m_d)$  and a variance vector  $\boldsymbol{\sigma} = (s_1, ..., s_d)$ , where each  $s_i \ge 0$ . The variational loss function for this output is given by

 $\frac{1}{2} \sum_{i=1,..d} \{ (s_i)^2 + (m_i)^2 - \ln[(s_i)^2] - 1 \}.$ 

(This equation is somewhat different from the book.) Show that this variation loss is minimized when  $\mu=0$  and  $\sigma=1$  (i.e. all the means are 0 and all variances are 1).

Question on Generative Adversarial Networks

Considering training a GAN with generator G(z) and discriminator D(G(z)). The perfect discriminator outputs 1 on a real instance, 0 on a fake instance. Figure 1 shows the training losses for two generator loss functions. In detail, for the first m generated points, i = 1,...,m. Each point in the plot shows the value of  $D(G(z_i))$ ,  $J_1(G)$  and  $J_2(G)$ , defined as follows.

- 1.  $J_1(G) = -1/m \sum_{i=1,...m} ln(D(G(z_i)))$ . Shown as the blue curve.
- 2.  $J_2(G) = 1/m \sum_{i=1,...,m} ln(1-D(G(z_i)))$ . Shown as the orange curve.



- 1. Early in the training, is the value of D(G(z)) closer to 0 or closer to 1? Explain why.
- 2. Which of the two cost functions would you choose to train your GAN? Justify your answer.
- 3. A GAN is successfully trained when D(G(z)) is close to 1. True or False? Explain your answer.

## Question on Seq2Seq Problems

Consider multiple-length seq2seq for an French-to-English MT program, as described in the text. We have decided on two sentence sizes

- Up to 8 words (including STOPs) for English, 10 for French
- Up to 10 words for English, 13 for French.

Suppose that the French input sentence is "A B C D E F" and the English output sentence is "M N O P Q R S T". Do you need to add one or more instances of a <STOP> token to make the input and output sentences formatted for processing by a neural network? If so, show

- the French input with <STOP> tokens added
- the English output with <STOP> tokens added

#### Question on Reinforcement Learning

Figure 6.1 in the text shows the value iteration algorithm. Show that the following *V* table gives the correct state values (to two significant digits) after the second pass of value iteration.

0	0	0	0
0	0	0	0
0	0	.1	0
0	.1	.43	0

[Figure 6.3 in the book has .46 in place of .43 for V(14). This correct if the states are updated sequentially in order and states 10 and 13 are updated before state 14. You can answer the question by replacing the V table above by Figure 6.3. Please specify which table you are using]

Answer Format: You can use the reasoning in the book and fill in further details. For example, the book explains why a state receives value 0 if all its neighbours currently have value 0. A nice way to answer this question is to show both Q-values and V-values. You could also show current\_state-action-next\_state triples for transitions where either the rewards or the next V values are non-zero.