Natural Language Understanding, Generation, and Machine Translation (2023–24)

2024 edition: Alexandra Birch, Jonas Waldendorf

Past editions: Christos Baziotis, Arturo Oncevay, Adam Lopez, Mirella Lapata, and

Frank Keller

School of Informatics, University of Edinburgh

NLU+ Coursework 1: Recurrent Neural Networks

This assignment is due on Friday, the 16th of February 2024, at 12pm, UTC.

Deadline for choosing partners: Tuesday, the 23rd of January 2024 at 5pm, UTC.

Executive Summary. This first part of this coursework will walk you through the implementation of a few critical parts of a recurrent neural network and the backpropagation algorithm, which are foundational components of the modern NLP toolkit. You will not need to derive anything mathematically—the maths will be given to you. This is similar to a lot of practical work in research and industry—it is *much* more common to implement and test an existing model from a specification than to derive a new one from first principles. In fact, what we ask you to implement is relatively low level, compared to what you can do with modern machine learning libraries. Nevertheless, the implementation-from-specification pattern is similar, and we hope that revealing more of the mathematical details will help demystify the underlying algorithms for you, and give you some degree of comfort with them. If you are very interested in deriving algorithms like this, we encourage you to dig further into the mathematics, but nothing in the coursework requires this.

The next parts of the coursework ask you to experiment with training regimes, and to adapt the model to an interesting psycholinguistic task that tests the model's behaviour on a phenomenon that humans process effortlessly—number agreement between subject and predicate in English. We strongly advise that you read the accompanying research paper from which the coursework dataset orginates (Linzen *et al.*, 2016). Using this setup you will compare the ability of recurrent neural networks and gate recurrent units to learn long range dependencies.

Finally, there is an open-ended final question for especially keen students. You do not need to answer the final question to receive a good mark, and should only attempt

it if you are confident in your solutions to the rest of the coursework and you have sufficient time remaining.

Submission Deadline and Pacing. The coursework is due on Friday, the 16th of February at 12pm (UTC), in week 5. If you want to work with a partner, you have an earlier deadline to relate this to us, on Tuesday, the 23rd of January at 12pm (UTC) i.e. three weeks before the coursework due date.

There are five questions. Simply answering the first four will earn you a good mark, and they can be easily done in the available time. This is an empirical observation from previous offerings of this course, which included a similar (though not identical!) coursework. The optional open-ended part of the coursework is Question 5. If you want to attempt it, you will need to finish the first four questions a week early. This is a more ambitious schedule, so plan accordingly.

Pair work policy. You are *strongly* encouraged to work with a partner on this assignment.¹ When you work with another person, you learn more, because you need to explain things to each other as you pool your collective expertise to solve problems. Explaining something to another person helps you debug your own thinking, and their questions help you overcome your own blind spots—something you cannot do on your own by staring at maths, code, or data. For this reason, it is best to seek partners with complementary skills to your own. You may not work in teams of three or more.

If you work with a partner, only one of you should submit your completed work. But I need to know that the submission represents the work of two people, and I need to know that reliably in advance in order to estimate marking hours. So, **you must do the following to ensure that both of you receive credit:**

- 1. Find a partner and confirm which one of you will submit the partnering form. If this is you, **confirm you have correctly recorded your partner's student number (UUN)**.
- 2. Submit your names and UUNs using the following form:

Link to form

This form is only accessible when you are logged into your University email account, so you may have to enter your password.

With this form, you can also notify us whether you want to be randomly paired with a partner.

You must do this by **Tuesday, the 23rd of January at 5pm (UTC)**. You may not change partners after you have done this, so take this commitment seriously. If I do not receive a

¹You will follow a similar partnering procedure for Coursework 2.

submission from you or your partner by the deadline, I will assume that you are working alone. Either way, I will confirm your choice with you shortly after the deadline to avoid later confusion. **I advise you to choose your partner now and get to work.** Piazza offers a feature that allows you to search for a partner, and you are welcome to use this. We can also assign you a partner if you need one. You can select that option in the partnering form.

Submission. Your solution should be delivered in two parts and uploaded to Gradescope. Learn. **Do not include your name or your partner's name in either the code or the writeup.** The coursework will be marked anonymously since this has been empirically shown to reduce bias. (I anonymize the filenames before marking.)

For your writeup:

- Write up your answers in a file titled <UUN>.pdf. For example, if your UUN is \$123456, your corresponding PDF should be named \$123456.pdf.
- The answers should be clearly numbered and can contain text, diagrams, graphs, formulas, as appropriate. Do not repeat the question text. If you are not comfortable with writing math on Latex/Word you are allowed to include scanned handwritten answers in your submitted pdf. You will lose marks if your handwritten answers are illegible.
- On Learn Ultra, under Assessment, access the Gradescope LTI link and select the assignment "Coursework 1 REPORT". Upload your <UUN>.pdf to this assignment, and use Gradescope to annotate the submission with the relevant question number for each page. To make this easier each page should contain the answer to only one question.
- Please make sure you have submitted the right file. We cannot make concessions for students who turn in incomplete or incorrect files by accident.

For your code and parameter files:

- On Learn Ultra, under Assessment, access the Gradescope LTI link and select the assignment "Coursework 1 CODE". Upload the rnn.py, gru.py and runner.py as well as your saved parameters rnn.U.npy, rnn.V.npy, and rnn.W.npy to the assignment.
- Please note that these Questions are marked using an autograder which requires the uploaded files to be in root the directory (either upload the files individually or create a zip of the files, do not place the files in a folder and then upload a zip of the folder). A test is provided that will tell you if the files have been uploaded with the correct file structure.

Good Scholarly Practice Please remember the University requirement as regards all assessed work for credit. Details and advice about this can be found at:

http://web.inf.ed.ac.uk/infweb/admin/policies/academic-misconduct

and links from there. Note that you are required to take reasonable measures to protect your assessed work from unauthorised access. For example, if you put your work in a public repository then you must restrict access only to yourself and your partner. For your writeup, and particularly on the final question, you should pay close attention to the guidance on plagiarism. In short: the litmus test for plagiarism is not the Turnitin check—that is simply an automated assistant. If you have borrowed or lightly edited someone else's words, you have plagiarised. Write your report in your own words. We are not marking for eloquence—as long as we can clearly understand what you did, that is fine.

School policy requires you to complete coursework yourself, using your own words, code, figures, etc. and to acknowledge any sources of text, code, figures etc. that are not your own. This policy includes the use of intelligent writing assistants such as ChatGPT. Using such an assistant without acknowledgement is a form of academic misconduct.

Assignment Data Files for this assignment are available on the course page in Learn.

Python Virtual Environment For this assignment you will be using Python along with a few open-source packages. These packages cannot be installed directly, so you will have to create a virtual environment. We are using virtual environments to make the installation of packages and retention of correct versions as simple as possible. For this assignment we are going to use Miniconda + Python 3.7. We will be using the same environment for both Assignment 1 and 2, so we recommend following these instructions to set yourself up for this course.

The instructions below are for DICE. You are free to use your own machine, but we cannot offer support for non-DICE machines. Similarly, this installation is for CPUs and we cannot offer support for GPU programming or any attempts to execute the assignment on a cluster such as mlp.inf.ed.ac.uk. The package specification file (requirements.txt) will only work for Linux based systems, if you choose to use a different machine then see the note below.

Open a terminal on a DICE machine and follow these instructions. We are expecting you to enter these commands in one-by-one. Waiting for each command to complete will help catch any unexpected warnings and errors. The total installation is about 4.33GB, please ensure you have sufficient space using the freespace command on DICE.

First install Miniconda from the home directory of your DICE user space. You can skip this stage if you already have Miniconda installed.

```
$> wget https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh
$> bash ./Miniconda3-latest-Linux-x86_64.sh
$> rm ./Miniconda3-latest-Linux-x86_64.sh
$> source ~/.bashrc
```

Now your default Python version should be 3.7. Confirm with python3 -version. Having ensured you have a right version of Python, proceed to create a new environment called nlu. We will use the file requirements.txt included in the assignment code to generate this environment.

1. Create a new folder for the assignment:

```
$> mkdir nlu_cw1
$> cd nlu cw1
```

- 2. Copy all code and data from Learn into this folder and unzip the compressed file as required.
- 3. Create an environment using the provided requirements.txt file: conda create -n nlu --file ./requirements.txt -c pytorch
- 4. Activate the nlu virtual environment:

```
$> conda activate nlu
```

5. **Optional** Clean your workspace to free up space:

```
$> conda clean -all
```

You should now have all the required packages installed. You only need to create the virtual environment and perform the package installations (step 1-3) **once**. However, make sure you activate your virtual environment (step 4) **every time** you open a new terminal to work on your assignment. Remember to use the conda deactivate command to disable the virtual environment when you don't need it.

1. Activate the environment:

```
$> conda activate nlu
```

2. Deactivate the environment (if you want to work on something else):

```
$> conda deactivate nlu
```

Using a different machine The above installation is designed for usage on DICE, but will also work for other Linux machines. If you are using Mac or Windows, you will have difficulty installing the required packages using the requirements.txt file

as they are not designed for your computer. You can replace commands 3 & 4 using the following commands:

- 1. \$> conda create -n nlu python=3.7.5
- 2. \$> conda activate nlu
- 3. \$> conda install tqdm=4.40.2 numpy=1.18.5 gensim=3.8.0 pandas=0.25.3 seaborn=0.9.0 matplotlib=3.1.1
- 4. \$> conda install pytorch=1.3.1 torchvision=0.4.2 cpuonly=1.0 -c pytorch

Terminology/definitions Most neural network components, as well as their architecture and functionality, can be described using *matrix* and *vector* mathematical operations. Matrix and vector notation in the literature is inconsistent, so for this assignment we will use the following conventions:

- (1) Matrices are assigned bold capital letters, e.g., U, V, W.
- (2) Vectors are written in bold lower-cased letters, e.g., \mathbf{x} , \mathbf{s} , \mathbf{net}_{in} , \mathbf{net}_{out} .
- (3) For a matrix **M** and a vector **v**, **Mv** represents their matrix-vector (dot) product.²
- (4) For vectors \mathbf{v} and \mathbf{w} of equal length n, $\mathbf{v} \circ \mathbf{w}$ represents their *element-wise product*:

$$\mathbf{v} \circ \mathbf{w} = [\mathbf{v}_0 \mathbf{w}_0, \mathbf{v}_1 \mathbf{w}_1, \dots, \mathbf{v}_n \mathbf{w}_n]$$

Similarly, $\mathbf{v} + \mathbf{w}$ and $\mathbf{v} - \mathbf{w}$ express element-wise addition and subtraction.

- (5) For vectors \mathbf{v} and \mathbf{w} , $\mathbf{v} \otimes \mathbf{w}$ represents their *outer product*.³
- (6) In recurrent neural networks, we process a *sequence* (or *time*), where each component is in a different state depending on the position in the sequence (or time step). We use the notation $\mathbf{M}^{(t)}$, $\mathbf{v}^{(t)}$ to refer to matrices and vectors at time step t.

Provided code and use of NumPy We provide a number of template files which you must use to write your code. rnn.py and gru.py are the implementation of the models, runner.py is used to train the models and define loss functions. We also provide an additional module, rnnmath.py, which consists of helper functions you can use. Finally, we provide test.py, which performs a very basic test of your code. Please familiarize yourself with the provided code and make sure you **don't** change the provided function signatures. Other classes should not be changed.

 $^{^2 \}verb|https://en.wikipedia.org/wiki/Matrix multiplication$

https://en.wikipedia.org/wiki/Outer product

Throughout this assignment, you are required to use NumPy methods for all matrix/vector operations. **Do not try to implement matrix/vector functionality on your own.** If you need help with NumPy, please refer to its documentation⁴ and look for answers on Google before asking on Piazza.

Introduction

For this assignment, you are asked to implement some basic functionality of a Recurrent Neural Network (RNN) for Language Modeling (LM). Given a word sequence w_1, w_2, \ldots, w_t , a language model predicts the next word w_{t+1} by modeling:

$$P(w_{t+1} | w_1, \ldots, w_t).$$

Below, you will be introduced to the main elements of a simple RNN for LM, based on the model proposed by Mikolov *et al.* (2010). In Question 2, you will implement its core word prediction functionality and its training by implementing a loss function and the model's gradient accumulation through backpropagation. In Question 3, you will adapt the model to the agreement prediction task as well as implementing a GRU (gate recurrent unit) cell Cho *et al.* (2014). Using this code you will compare the ability of RNNs and GRUs to learn long range dependencies.

Extra Resources:

- For understanding RNNs and BPTT please read Guo (2013). It is a fairly formal treatment but it explains everything in detail and defines all the notation.
- For understanding LSTMs I recommend Chapter 4 of Alex Graves thesis Graves (2012).
- This tutorial is also useful Chistopher (2015).

Recurrent Neural Networks

A recurrent neural network for language modeling uses feedback information in the hidden layer to model the "history" w_1, w_2, \ldots, w_t in order to predict w_{t+1} . Formally, at each time step, the model needs to compute:

⁴http://docs.scipy.org/doc/numpy/reference/index.html

$$\mathbf{s}^{(t)} = f\left(\mathbf{net}_{in}^{(t)}\right) \tag{1}$$

$$\mathbf{net}_{in}^{(t)} = \mathbf{V}\mathbf{x}^{(t)} + \mathbf{U}\mathbf{s}^{(t-1)}$$
 (2)

$$\mathbf{net}_{in}^{(t)} = \mathbf{V}\mathbf{x}^{(t)} + \mathbf{U}\mathbf{s}^{(t-1)}$$

$$\mathbf{\hat{y}}^{(t)} = g\left(\mathbf{net}_{out}^{(t)}\right)$$

$$\mathbf{net}_{out}^{(t)} = \mathbf{W}\mathbf{s}^{(t)}$$
(3)

$$\mathbf{net}_{out}^{(t)} = \mathbf{Ws}^{(t)} \tag{4}$$

where f() and g() are the *sigmoid* and *softmax* activation functions respectively, $\mathbf{x}^{(t)}$ is the one-hot vector representing the vocabulary index of the word w_t , $\mathbf{net}_{in}^{(t)}$ and $\mathbf{net}_{out}^{(t)}$ are the activations for the hidden and output layers, and $\mathbf{s}^{(t)}$ and $\mathbf{\hat{v}}^{(t)}$ are the corresponding hidden and output vectors produced after applying the sigmoid and softmax nonlinearities.

For a given input $[w_1, w_2, \dots, w_t]$, the probability of the next word at time step t+1 can be read from the output vector $\hat{\mathbf{y}}^{(t)}$:

$$P(w_{t+1} = j \mid w_t, \dots, w_1) = \hat{y}_j^{(t)}$$
(5)

The parameters to be learned are:

$$\mathbf{U} \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{V} \in \mathbb{R}^{D_h \times |V|} \quad \mathbf{W} \in \mathbb{R}^{|V| \times D_h}$$
 (6)

where **U** is the matrix for the recurrent hidden layer, **V** is the input word representation matrix, W is the output word representation matrix, and D_h is the dimensionality of the hidden layer.

Question 1: Training RNNs [20 marks]

When training RNNs, we need to propagate the errors observed at the output layer \hat{y} back through the network, and adjust the weight matrices U, V and W to minimize the observed loss w.r.t. a desired output. There are several loss functions suitable for use in RNNs. In RNN language models, an effective loss function is the cross-entropy loss:

$$J^{(t)}(\theta) = -\sum_{i=1}^{|V|} d_j^{(t)} \log \hat{y}_j^{(t)}$$
 (7)

where $\mathbf{d}^{(t)} = [d_1^{(t)}, d_2^{(t)}, \dots, d_{|V|}^{(t)}]$ is the one-hot vector representing the vocabulary index of desired output word at time t. In order to evaluate the model's performance, we average the cross-entropy loss across all steps in a sentence and across all sentences in the dataset.

- (a) In the file rnn.py, implement the method predict of the RNN class. The method is used for *forward prediction* in your RNN and takes as input a sentence as a list of word indices $[w_1, \dots, w_n]$. The return values are the matrices produced by concatenating hidden vectors $\mathbf{s}^{(t)}$ and output vectors $\hat{\mathbf{y}}^{(t)}$ for $t = 1, 2, \dots, n$. [5 marks]
- (b) In runner.py, implement the methods compute_loss and compute_mean_loss. Given a sequence of input words $w = [w_1, ..., w_n]$ and a sequence of desired output words $d = [d_1, ..., d_n]$, compute_loss should return the total loss produced by the model's predictions for the sentence. The compute_mean_loss should compute the average loss over a corpus of input sentences. It should average across all words in all sentences of the given corpus. [5 marks]

Optimizing the loss using backpropagation means we have to calculate the update values Δ w.r.t. the gradients of our loss function for the observed errors. For the output layer weights, at time step t we accumulate the matrix **W** updates using:

$$\Delta \mathbf{W} = \eta \sum_{p=1}^{n} \delta_{out,p}^{(t)} \otimes \mathbf{s}_{p}^{(t)}$$
 (8)

$$\boldsymbol{\delta}_{out,p}^{(t)} = (\mathbf{d}_p^{(t)} - \mathbf{\hat{y}}_p^{(t)}) \circ g'(\mathbf{net}_{out,p}^{(t)})$$

$$\tag{9}$$

where η is the learning rate and p indicates the index of the current training pattern (sentence). We then further propagate the error observed at the output back to V with:

$$\Delta \mathbf{V} = \eta \sum_{n=1}^{n} \delta_{in,p}^{(t)} \otimes \mathbf{x}_{p}^{(t)}$$
 (10)

$$\delta_{in,p}^{(t)} = \mathbf{W}^T \delta_{out,p}^{(t)} \circ f'(\mathbf{net}_{in,p}^{(t)})$$
(11)

The derivatives of the softmax and sigmoid functions are respectively given as⁶:

$$g'(\mathbf{net}_{out,p}^{(t)}) = \vec{1} \tag{12}$$

$$f'(\mathbf{net}_{in,p}^{(t)}) = \mathbf{s}_p^{(t)} \circ (\vec{1} - \mathbf{s}_p^{(t)})$$
 (13)

⁵See the provided documentation on rnn.py for more details on the functions you need to implement.

 $^{^{6}}$ We use $\vec{1}$ as shorthand for the all-ones vector of appropriate length.

Finally, in order to update the recurrent weights U, we need to look back one step in time:

$$\Delta \mathbf{U} = \eta \sum_{p=1}^{n} \delta_{in,p}^{(t)} \otimes \mathbf{s}_{p}^{(t-1)}$$
 (14)

(c) In rnn.py, implement the method acc_deltas that accumulates the weight updates for **U**, **V** and **W** for a *truncated* backpropagation through the RNN, where we only look back one step in time as described above, rather than the entire computation graph. [5 marks]

Now we have implemented truncated backpropagation (BP) for recurrent networks—that is, RNNs that just look at the previous hidden layer when accumulating ΔU and ΔV . We can extend truncated backpropagation to look at the previous τ time steps during backpropagation. At time t, the updates ΔW can be derived as before. For ΔU and ΔV , we additionally recursively update at times (t-1), (t-2) \cdots $(t-\tau)$:

$$\Delta \mathbf{V} = \eta \sum_{p=1}^{n} \delta_{in,p}^{(t-\tau)} \otimes \mathbf{x}_{p}^{(t-\tau)}$$
(15)

$$\Delta \mathbf{U} = \eta \sum_{p=1}^{n} \delta_{in,p}^{(t-\tau)} \otimes \mathbf{s}_{p}^{(t-\tau-1)}$$
(16)

$$\delta_{in,p}^{(t-\tau)} = \mathbf{U}^T \delta_{in,p}^{(t-\tau+1)} \circ f'(\mathbf{net}_{in,p}^{(t-\tau)})$$
(17)

(d) Implement the method acc_deltas_bptt that accumulates the weight updates for U, V and W using backpropagation through time for τ time steps. [5 marks]

There's one last thing you'll need to do before your models will train: you'll need to complete the implementation of the train-lm-rnn mode in the __main__ method of the code in runner.py. For this, you will minimally need to instantiate the RNN class, instatiate the Runner class with the RNN and call the train method on the RNN with the appropriate arguments, and save the resulting matrices. You may find it useful for your own understanding to log different aspects of the training process here.

You do not need to report anything in your writeup for Question 1. It will be evaluated solely on the basis of your code.

Question 2: Language Modeling [15 marks]

By now you should have everything in place to train a full Recurrent Neural Network using backpropagation through time. In the following questions, we will use the training

and development data provided in wiki-train.txt and wiki-dev.txt. The training data consists of sentences from the parsed English Wikipedia corpus from Linzen *et al.* (2016), and each input/output pair x, d is of the form $[w_1, \dots w_n] / [w_2, \dots w_{n+1}]$ that is, the desired output is always the next word of the current input:

time index	t=1	t=2	t=3	t=4
input:	Banks	struggled	with	the
output:	struggled	with	the	crisis

The utils.py module provides functions to read the Wikipedia data, and the __main__ method of the runner.py module provides some starter code for training your models. Use mode train-lm-rnn to train your language model. In order to start the training code you first have to instantiate the RNN and then a runner in order to call the train method.

(a) Perform parameter tuning using a subset of the training and development sets. You must use a fixed vocabulary of size 2,000, and consider all combinations of the following parameter settings:

learning rate: 0.5, 0.1, or 0.05 number of hidden units: 25 or 50 number of steps to look back in truncated backpropagation: 0, 2, or 5

The mode train-lm-rnn in runner.py allows for more parameters, which you are free to explore. You should tune your model to maximize generalization performance (minimize cross-entropy loss). For these experiments, use the first 1,000 sentences of both the training and development sets and train for 10 epochs.⁷ Report your findings and interpret them. Your interpretation need not be more than a paragraph. [10 marks]

(b) Using your best parameter settings found in (a), train an RNN on a much larger training set. Use a fixed vocabulary size of 2000, train on 25,000 sentences, and, as before, use the first 1,000 development sentences to evaluate the model's performance during training. When your model is trained⁸, evaluate it on the **test** set and report the mean loss, as well as both the adjusted and unadjusted perplexity your model achieves⁹. Save your final learned matrices **U**, **V** and **W** as files rnn.U.npy, rnn.V.npy and rnn.W.npy, respectively. [5 marks]

⁷Note that training models might take some time. For example, a sweep of the parameters settings described above should take roughly 2 hours on a student lab DICE machine. Please avoid using student.compute to train your models as run times will become very slow on a busy server.

⁸This should also take roughly 2 hours. Again, avoid using student.compute.

⁹Use your method compute_mean_loss to calculate loss on the development set, and the provided method adjust_loss to get adjusted/unadjusted perplexities for your models

Question 3: Predicting Agreement with RNNs and GRUs [15 marks]

The form of an English third-person present tense <u>verb</u> depends on whether the **head** of the syntactic subject is plural or singular. For example, native English speakers strongly prefer sentences (i) and (iv) below, and regard (ii) and (iii) as ungrammatical, as indicated by the *:

- i) The **key** is on the table.
- ii) *The **key** <u>are</u> on the table.
- iii) *The **keys** is on the table.
- iv) The keys are on the table.

This agreement tends to persist even when the head of the subject is separated from the verb by intervening words:

v) The **keys** to the cabinet are on the table.

Agreement rules like this occur in many languages, and are often more complex than in English. Our goal for this question will be to test (in a limited way) whether an RNN can learn them. For our first test, we will train a model predict agreement using direct supervision. That is, we will give our model the sequence of words preceding the verb, and we will ask it to predict whether the verb is singular (VBZ), or plural (VBP). Our training and test data will be in this form:

time index	t=1	t=2	t=3	t=4	t=5	
input:	The	keys	to	the	cabinet	
output:					VBP	

Since the task is now binary classification, we must make some changes to the RNN. Instead of making predictions at **every** time step, we only make a prediction at the **final** time step.

(a) Implement new functions for weight updates (acc_deltas_np, acc_deltas_bptt_np, rnn.py), loss fuction (compute_loss_np, runner.py), and prediction accuracy (compute_acc_np, runner.py) to reflect the structure of the number prediction problem. [5 marks]

Since the head of the subject may be arbitrarily far from the verb, this problem is a natural application of RNNs. The model needs to learn long-term dependencies to correctly predict whether a verb is singular or plural.

Gated Recurrent Unit

The Gated Recurrent Unit (GRU) is special type of RNN which is capable of learning long range dependencies. The overall structure is the same but each RNN cell is replaced by a GRU cell. Formally, at each time step the model computes:

$$\mathbf{r}^{(t)} = f\left(\mathbf{V}_r \mathbf{x}^{(t)} + \mathbf{U}_r \mathbf{s}^{(t-1)}\right)$$
(18)

$$\mathbf{z}^{(t)} = f\left(\mathbf{V}_z\mathbf{x}^{(t)} + \mathbf{U}_z\mathbf{s}^{(t-1)}\right)$$
 (19)

$$\tilde{\mathbf{h}}^{(t)} = tanh\left(\mathbf{V}_h\mathbf{x}^{(t)} + \mathbf{U}_h\left(\mathbf{r}^{(t)} \circ \mathbf{s}^{(t-1)}\right)\right)$$
(20)

$$\mathbf{s}^{(t)} = \mathbf{z}^{(t)} \circ \mathbf{s}^{(t-1)} + \left(1 - \mathbf{z}^{(t)}\right) \circ \tilde{\mathbf{h}}^{(t)}$$
(21)

$$\mathbf{net}_{out}^{(t)} = \mathbf{Ws}^{(t)} \tag{22}$$

$$\hat{\mathbf{y}}^{(t)} = g\left(\mathbf{net}_{out}^{(t)}\right) \tag{23}$$

where f(), g() and tanh() are the sigmoid, softmax and tanh activation functions respectively. $\mathbf{x}^{(t)}$ is the one-hot vector representing the vocabulary index of the word w_t . $\mathbf{r}_{in}^{(t)}$ and $\mathbf{z}_{in}^{(t)}$ are the reset and update gates respectively. $\tilde{\mathbf{h}}^{(t)}$ is an additional "candidate" hidden state that uses the reset gate to remove irrelevant information. $\mathbf{s}^{(t)}$ is the actual hidden state which is obtained using the update gate and $\tilde{\mathbf{h}}^{(t)}$. $\mathbf{net}_{out}^{(t)}$ is the activation of the output layer.

The parameters to be learned are:

$$\mathbf{U}_r \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{V}_r \in \mathbb{R}^{D_h \times |V|} \tag{24}$$

$$\mathbf{U}_z \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{V}_z \in \mathbb{R}^{D_h \times |V|} \tag{25}$$

$$\mathbf{U}_{r} \in \mathbb{R}^{D_{h} \times D_{h}} \quad \mathbf{V}_{r} \in \mathbb{R}^{D_{h} \times |V|}$$

$$\mathbf{U}_{z} \in \mathbb{R}^{D_{h} \times D_{h}} \quad \mathbf{V}_{z} \in \mathbb{R}^{D_{h} \times |V|}$$

$$\mathbf{U}_{h} \in \mathbb{R}^{D_{h} \times D_{h}} \quad \mathbf{V}_{h} \in \mathbb{R}^{D_{h} \times |V|}$$

$$(24)$$

$$(25)$$

$$(26)$$

$$\mathbf{W} \in \mathbb{R}^{|V| \times D_h} \tag{27}$$

where U is the matrix for the recurrent hidden layer, V is the input word representation matrix, W is the output word representation matrix, and D_h is the dimensionality of the hidden layer. Unlike the RNN, this implementation of a GRU learns three separate sets of U and V parameters denoted by r, z and h. These subscripts represent the reset, update and hidden state parameters respectively.

(b) In the file gru.py, implement the forward method. Note this method only needs to perform the forward prediction for a single time step, the prediction for an entire sequence is already handled by gru_abstract.py if the forward method is correctly implemented. It takes as input the single word index for the current time step w_t and the hidden state vector from the previous time step $\mathbf{s}^{(t-1)}$. The return values are the output vector $\mathbf{y}^{(t)}$, the reset gate vector $\mathbf{r}^{(t)}$, the update gate vector $\mathbf{z}^{(t)}$, the "candidate" hidden state vector $\tilde{\mathbf{h}}^{(t)}$ and the new hidden state vector $\mathbf{s}^{(t)}$.

Next, implement the acc_deltas_np and acc_deltas_bptt_np methods in the gru.py class. The gru.py class provides a self.backward() method that takes $\delta_{out,p}^{(t)}$ given in Equation 9 as an input and performs the backpropagation for the rest of the GRU model. Specifically, self.backward() takes as input the list of word indicies x, the current time step t, the gradients of the output $\delta_{out,p}^{(t)}$ and an optional parameter steps. steps defines the number of steps for BPTT and is set to 0 by default. As for the RNN acc_deltas_np and acc_deltas_bptt_np should reflect the structure of the number prediction problem. [5 marks]

You'll need to complete the implementation of the train-np-rnn and train-np-gru modes in the __main__ method of the runner.py file. Check your implementation of the RNN by running the code:

\$> python runner.py train-np-rnn data_dir hdim lookback learning_rate
and the GRU by running:

- \$> python runner.py train-np-gru data_dir hdim lookback learning_rate
- (c) Using the first 10,000 training sentences, and a learning rate of 0.5 training, train your models for 10 epochs and compare the accuracy of the RNN and the GRU models by performing a sweep of the number of hidden units, setting them to: 10, 25 and 50.

Report your the results for both models and briefly describe how they differ from each other. [5 marks]

Question 4: Comparing Recurrent Models [10 marks]

In Question 3, you tested GRUs and RNNs on the agreement task. However, you did not use backpropagation through time (the number of BPTT steps was set to 0). The aim of this question is to compare the two models when backpropagation through time is used. Specifically, you should compare how the two architectures perform on the number prediction task as number of BPTT steps increases.

Using the first 10,000 training sentences, a learning rate of 0.5 training and 50 hidden units, train your models for 10 epochs and compare the RNN and the GRU models whilst varying the BPTT steps. You may need to change the runner class to log additional

information during training; in this case we recommend doing this in a new runner class to ensure that any changes to the code do not break the automatic marking.

Report your results, and provide **an interpretation** of the results based on the differences between the two architectures, this should be no more than two paragraphs. This question is more open-ended than the previous questions and very good answers should consider the advantages of GRU cells, how the task can be used to demonstrate them and what results are needed to support your argument.

Please note that the code will not be marked for this question; we will mark only the writeup.

Question 5: Exploration [25 marks]

The final part of this question is open-ended, and it is intended primarily for students who want to deepen their knowledge, for the price of some additional work. **You are not required to attempt this part of the assignment**, and you will receive a good mark if you do a good job on only the first four questions. If you are not yet comfortable with the material in Questions 1–4, focus on those questions only.

I want you to ask your **own** question about the number prediction task and RNNs, and then attempt to answer it. Put another way: I want you to develop a simple scientific hypothesis and test it. **The final writeup should advance a claim about the answer, supported by empirical evidence from your experiements**. You must analyze this evidence, so you'll need to look at the data, you'll probably need to measure one or more properties of interest, and you'll need to interpret the results with respect to your stated question. You need to reason scientifically here, and **design experiments that test your hypothesis**. For example, if you hypothesize that LSTMs "solve the vanishing gradient problem" (something that many, *many* students have claimed on past versions of this assignment, with little insight), then your evidence had better include an analysis of the gradients. Better test perplexity is *not* evidence that gradients are propagating over longer distances. ¹⁰

You are welcome (but absolutely not required) to implement another model, if doing so helps you answer a question. But this question is not a test your implementation skills, and you will not receive more marks simply for writing many lines of code, and indeed the best answers to this question have typically required no code at all. A simple, well-posed analysis or experiment, clearly explained, is much better than an implementation of a fancy architecture.

¹⁰The hypothesis doesn't have to be original—reproducing scientific results is important! But you should not expect to earn many marks with a poor analysis of a well-established hypothesiss.

This question is also not a test of the *amount* of work that you do: I realize that you have limited time for this part of the coursework, and I don't expect you to produce a novel result or even a novel question. I *do* expect you to ask a small question that you find interesting, even if it's a question you've seen elsewhere, or just something small that intrigued you about the data. It is perfectly ok if you present an interesting question and run a well-posed analysis whose results are inconclusive. That is how most science works. In short, the brief for this question is for you to be curious and engage with the topic of the course. Your writeup should tell us something interesting that you learned (or attempted to learn).

This question will be marked according to the descriptors for the common marking scheme. Note that these criteria require "elements of personal insight / creativity / originality" for a mark above 70—that's why this question is worth 25 marks. For the same reason, it will be rare for answers to this question to receive more than 15 points. That's only 2% of your overall mark in the course, so **I urge you not spend too much time on it**. Your answer to this question should not be more than two pages, including figures. [25 marks]

References

Olah Chistopher. Understanding LSTM Networks. https://colah.github.io/posts/2015-08-Understanding-LSTMs/, 2015. Unpubl. Blog.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics.

Alex Graves. Supervised Sequence Labelling with Recurrent Neural Networks. PhD thesis, Technische Universität Munchen, 2012.

Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. Colorless green recurrent networks dream hierarchically. In *Proceedings of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018.

¹¹"Research is the process of going up alleys to see if they are blind."—Marston Bates

¹²https://web.inf.ed.ac.uk/infweb/student-services/taught-students/
information-for-students/information-for-all-students/your-studies/
common-marking-scheme

Jiang Guo. Backpropagation Through Time. http://ir.hit.edu.cn/~jguo/docs/notes/bptt.pdf, 2013. Unpubl. ms., Harbin Institute of Technology.

Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics*, 4:521–535, 2016.

Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernockỳ, and Sanjeev Khudanpur. Recurrent neural network based language model. In *INTERSPEECH*, volume 2, page 3, 2010.

Acknowledgements

This assignment is based in parts on code and text kindly provided by Richard Socher, from the course on "Deep Learning for Natural Language Processing". 13

¹³http://cs224d.stanford.edu/