

# CS 224n: Assignment #2

**Due date: 2/8 11:59 PM PST** (You are allowed to use 3 late days maximum for this assignment)

These questions require thought, but do not require long answers. Please be as concise as possible.

? We encourage students to discuss in groups for assignments. **However, each student must finish the problem set and programming assignment individually, and must turn in her/his assignment.** We ask that you abide the university Honor Code and that of the Computer Science department, and make sure that all of your submitted work are done by yourself.

Please review any additional instructions posted on the assignment page at <http://cs224n.stanford.edu/assignments.html>. When you are ready to submit, please follow the instructions on Piazza.

## 1 Tensorflow Softmax (25 points)

In this question, we will implement a linear classifier with loss function

$$J(\mathbf{W}) = CE(\mathbf{y}, \text{softmax}(\mathbf{x}\mathbf{W} + \mathbf{b}))$$

Where  $\mathbf{x}$  is a vector of features,  $\mathbf{W}$  is the model's weight matrix, and  $\mathbf{b}$  is a bias term. We will use TensorFlow's automatic differentiation capability to fit this model to provided data.

- (a) (5 points, coding) Implement the softmax function using TensorFlow in `q1_softmax.py`. Remember that

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Note that you may **not** use `tf.nn.softmax` or related built-in functions. You can run basic (non-exhaustive tests) by running `python q1_softmax.py`.

- (b) (5 points, coding) Implement the cross-entropy loss using TensorFlow in `q1_softmax.py`. Remember that

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{N_c} y_i \log(\hat{y}_i)$$

where  $\mathbf{y} \in \mathbb{R}^{N_c}$  is a one-hot label vector and  $N_c$  is the number of classes. This loss is summed over all examples in a minibatch. Note that you may **not** use TensorFlow's built-in cross-entropy functions for this question. You can run basic (non-exhaustive tests) by running `python q1_softmax.py`.

- (c) (5 points, coding/written) Carefully study the Model class in `model.py`. Briefly explain the purpose of placeholder variables and feed dictionaries in TensorFlow computations. Fill in the implementations for `add_placeholders` and `create_feed_dict` in `q1_classifier.py`.

**Hint:** Note that configuration variables are stored in the Config class. You will need to use these configuration variables in the code.

(d) (5 points, coding) Implement the transformation for a softmax classifier in the function `add_prediction_op` in `q1_classifier.py`. Add cross-entropy loss in the function `add_loss_op` in the same file. Use the implementations from the earlier parts of the problem (already imported for you), **not** TensorFlow built-ins.

(e) (5 points, coding/written) Fill in the implementation for `add_training_op` in `q1_classifier.py`. Explain in a few sentences what happens when the model's `train_op` is called (what gets computed during forward propagation, what gets computed during backpropagation, and what will have changed after the op has been run?). Verify that your model is able to fit to synthetic data by running `python q1_classifier.py` and making sure that the tests pass.

**Hint:** Make sure to use the learning rate specified in `Config`.

## 2 Neural Transition-Based Dependency Parsing (50 points)

In this section, you'll be implementing a neural-network based dependency parser. A dependency parser analyzes the grammatical structure of a sentence, establishing relationships between "head" words and words which modify those heads. Your implementation will be **a transition-based parser, which incrementally builds up a parse one step at a time**. At every step it maintains a partial parse, which is represented as follows:

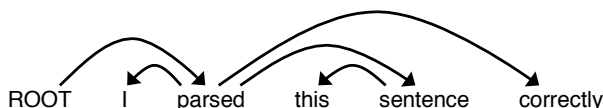
- A *stack* of words that are currently being processed.
- A *buffer* of words yet to be processed.
- A list of *dependencies* predicted by the parser.

Initially, the stack only contains ROOT, the dependencies lists is empty, and the buffer contains all words of the sentence in order. At each step, the parse applies a *transition* to the partial parse until its buffer is empty and the stack is size 1. The following transitions can be applied:

- **SHIFT:** removes the first word from the buffer and pushes it onto the stack.
- **LEFT-ARC:** marks the second (second most recently added) item on the stack as a dependent of the first item and removes the second item from the stack.
- **RIGHT-ARC:** marks the first (most recently added) item on the stack as a dependent of the second item and removes the first item from the stack.

Your parser will decide among transitions at each state using a neural network classifier. First, you will implement the partial parse representation and transition functions.

(a) (6 points, written) Go through the sequence of transitions needed for parsing the sentence "*I parsed this sentence correctly*". The dependency tree for the sentence is shown below. At each step, give the configuration of the stack and buffer, as well as what transition was applied this step and what new dependency was added (if any). The first three steps are provided below as an example.



stack	buffer	new dependency	transition
[ROOT]	[I, parsed, this, sentence, correctly]		Initial Configuration
[ROOT, I]	[parsed, this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]		SHIFT
[ROOT, parsed]	[this, sentence, correctly]	parsed→I	LEFT-ARC

- (b) (2 points, written) A sentence containing  $n$  words will be parsed in how many steps (in terms of  $n$ )? Briefly explain why.
- (c) (6 points, coding) Implement the `__init__` and `parse_step` functions in the `PartialParse` class in `q2_parser_transitions.py`. This implements the transition mechanics your parser will use. You can run basic (not-exhaustive) tests by running `python q2_parser_transitions.py`.
- (d) (6 points, coding) Our network will predict which transition should be applied next to a partial parse. We could use it to parse a single sentence by applying predicted transitions until the parse is complete. However, neural networks run much more efficiently when making predictions about batches of data at a time (i.e., predicting the next transition for a many different partial parses simultaneously). We can parse sentences in minibatches with the following algorithm.

---

**Algorithm 1** Minibatch Dependency Parsing

---

**Input:** `sentences`, a list of sentences to be parsed and `model`, our model that makes parse decisions

Initialize `partial_pares` as a list of partial parses, one for each sentence in `sentences`

Initialize `unfinished_pares` as a shallow copy of `partial_pares`

**while** `unfinished_pares` is not empty **do**

    Take the first `batch_size` parses in `unfinished_pares` as a minibatch

    Use the `model` to predict the next transition for each partial parse in the minibatch

    Perform a parse step on each partial parse in the minibatch with its predicted transition

    Remove the completed (empty buffer and stack of size 1) parses from `unfinished_pares`

**end while**

**Return:** The dependencies for each (now completed) parse in `partial_pares`.

---

Implement this algorithm in the `minibatch_parse` function in `q2_parser_transitions.py`. You can run basic (not-exhaustive) tests by running `python q2_parser_transitions.py`.

*Note: You will need `minibatch_parse` to be correctly implemented to evaluate the model you will build in part (h). However, you do not need it to train the model, so you should be able to complete most of part (h) even if `minibatch_parse` is not implemented yet.*

We are now going to train a neural network to predict, given the state of the stack, buffer, and dependencies, which transition should be applied next. First, the model extracts a feature vector representing the current state. We will be using the feature set presented in the original neural dependency parsing paper: *A Fast and Accurate Dependency Parser using Neural Networks*<sup>1</sup>. The function extracting these features has been implemented for you in `parser_utils`. This feature vector consists of a list of tokens (e.g., the last word in the stack, first word in the buffer, dependent of the second-to-last word in the stack if there is one, etc.). They can be represented as a list of integers

$$[w_1, w_2, \dots, w_m]$$

where  $m$  is the number of features and each  $0 \leq w_i < |V|$  is the index of a token in the vocabulary ( $|V|$  is the vocabulary size). First our network looks up an embedding for each word and concatenates them into a single input vector:

$$\mathbf{x} = [\mathbf{E}_{w_1}, \dots, \mathbf{E}_{w_m}] \in \mathbb{R}^{dm}$$

---

<sup>1</sup>Chen and Manning, 2014, <http://cs.stanford.edu/people/danqi/papers/emnlp2014.pdf>

where  $\mathbf{E} \in \mathbb{R}^{|V| \times d}$  is an embedding matrix with each row  $\mathbf{E}_i$  as the vector for a particular word  $i$ . We then compute our prediction as:

$$\begin{aligned}\mathbf{h} &= \text{ReLU}(\mathbf{x}\mathbf{W} + \mathbf{b}_1) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{h}\mathbf{U} + \mathbf{b}_2)\end{aligned}$$

(recall that  $\text{ReLU}(z) = \max(z, 0)$ ). We will train the model to minimize cross-entropy loss:

$$J(\theta) = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{N_c} y_i \log \hat{y}_i$$

To compute the loss for the training set, we average this  $J(\theta)$  across all training examples.

- (e) (4 points, coding) In order to avoid neurons becoming too correlated and ending up in poor local minima, it is often helpful to randomly initialize parameters. One of the most frequent initializations used is called Xavier initialization<sup>2</sup>.

Given a matrix  $\mathbf{A}$  of dimension  $m \times n$ , Xavier initialization selects values  $A_{ij}$  uniformly from  $[-\epsilon, \epsilon]$ , where

$$\epsilon = \frac{\sqrt{6}}{\sqrt{m+n}}$$

Implement the initialization in `xavier_weight_init` in `q2_initialization.py`. You can run basic (nonexhaustive tests) by running `python q2_initialization.py`. This function will be used to initialize  $\mathbf{W}$  and  $\mathbf{U}$ .

- (f) (2 points, written) We will regularize our network by applying Dropout<sup>3</sup>. During training this randomly sets units in the hidden layer  $\mathbf{h}$  to zero with probability  $p_{drop}$  and then multiplies  $\mathbf{h}$  by a constant  $\gamma$  (dropping different units each minibatch). We can write this as

$$\mathbf{h}_{drop} = \gamma \mathbf{d} \circ \mathbf{h}$$

where  $\mathbf{d} \in \{0, 1\}^{D_h}$  ( $D_h$  is the size of  $\mathbf{h}$ ) is a mask vector where each entry is 0 with probability  $p_{drop}$  and 1 with probability  $(1 - p_{drop})$ .  $\gamma$  is chosen such that the value of  $\mathbf{h}_{drop}$  in expectation equals  $\mathbf{h}$ :

$$\mathbb{E}_{p_{drop}}[\mathbf{h}_{drop}]_i = h_i$$

for all  $0 < i < D_h$ . What must  $\gamma$  equal in terms of  $p_{drop}$ ? Briefly justify your answer.

- (g) (4 points, written) We will train our model using the Adam<sup>4</sup> optimizer. Recall that standard SGD uses the update rule

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} J_{minibatch}(\boldsymbol{\theta})$$

where  $\boldsymbol{\theta}$  is a vector containing all of the model parameters,  $J$  is the loss function,  $\nabla_{\boldsymbol{\theta}} J_{minibatch}(\boldsymbol{\theta})$  is the gradient of the loss function with respect to the parameters on a minibatch of data, and  $\alpha$  is the learning rate. Adam uses a more sophisticated update rule with two additional steps<sup>5</sup>.

---

<sup>2</sup>This is also referred to as Glorot initialization and was initially described in <http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>

<sup>3</sup>Srivastava et al., 2014, <https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

<sup>4</sup>Kingma and Ma, 2015, <https://arxiv.org/pdf/1412.6980.pdf>

<sup>5</sup>The actual Adam update uses a few additional tricks that are less important, but we won't worry about them for this problem.

- (i) First, Adam uses a trick called *momentum* by keeping track of  $\mathbf{m}$ , a rolling average of the gradients:

$$\begin{aligned}\mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta) \\ \theta &\leftarrow \theta - \alpha \mathbf{m}\end{aligned}$$

where  $\beta_1$  is a hyperparameter between 0 and 1 (often set to 0.9). Briefly explain (you don't need to prove mathematically, just give an intuition) how using  $\mathbf{m}$  stops the updates from varying as much. Why might this help with learning?

- (ii) Adam also uses *adaptive learning rates* by keeping track of  $\mathbf{v}$ , a rolling average of the magnitudes of the gradients:

$$\begin{aligned}\mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta) \\ \mathbf{v} &\leftarrow \beta_2 \mathbf{v} + (1 - \beta_2) (\nabla_{\theta} J_{\text{minibatch}}(\theta) \circ \nabla_{\theta} J_{\text{minibatch}}(\theta)) \\ \theta &\leftarrow \theta - \alpha \circ \mathbf{m} / \sqrt{\mathbf{v}}\end{aligned}$$

where  $\circ$  and  $/$  denote elementwise multiplication and division (so  $\mathbf{z} \circ \mathbf{z}$  is elementwise squaring) and  $\beta_2$  is a hyperparameter between 0 and 1 (often set to 0.99). Since Adam divides the update by  $\sqrt{\mathbf{v}}$ , which of the model parameters will get larger updates? Why might this help with learning?

- (h) (20 points, coding/written) In `q2_parser_model.py` implement the neural network classifier governing the dependency parser by filling in the appropriate sections. We will train and evaluate our model on the Penn Treebank (annotated with Universal Dependencies). Run `python q2_parser_model.py` to train your model and compute predictions on the test data (make sure to turn off debug settings when doing final evaluation).

#### Hints:

- When debugging, pass the keyword argument `debug=True` to the main method (it is set to `true` by default). This will cause the code to run over a small subset of the data, so the training the model won't take as long.
- This code should run within 1 hour on a CPU.
- When running with `debug=True`, you should be able to get a loss smaller than 0.2 and a UAS larger than 65 on the dev set (although in rare cases your results may be lower, there is some randomness when training). When running with `debug=False`, you should be able to get a loss smaller than 0.08 on the train set and an Unlabeled Attachment Score larger than 87 on the dev set. For comparison, the model in the original neural dependency parsing paper gets 92.5 UAS. If you want, you can tweak the hyperparameters for your model (hidden layer size, hyperparameters for Adam, number of epochs, etc.) to improve the performance (but you are not required to do so).

#### Deliverables:

- Working implementation of the neural dependency parser in `q2_parser_model.py`. (We'll look at, and possibly run this code for grading).
  - Report the best UAS your model achieves on the dev set and the UAS it achieves on the test set.
  - List of predicted labels for the test set in the file `q2_test.predicted`.
- (i) **Bonus** (2 points). Add an extension to your model (e.g., l2 regularization, an additional hidden layer) and report the change in UAS on the dev set. Briefly explain what your extension is, what UAS the resulting model gets on the dev and test sets, and why the extension helps (or hurts!) the

model. Some extensions may require tweaking the hyperparameters in `Config` to make them effective. For both bonus points, your extended model should get better UAS than the baseline model without an extension. You should turn in your extension with your code, but **your extension should be turned off by default**. An easy way of doing this is add an `extension_on` boolean to your model's `Config` that is set to `False` when you turn in the code. If your extension does not improve over your baseline, do not turn in those predictions when you submit `q2_test.predicted`, use the baseline's ones instead (we will make sure these predictions produce a good UAS when grading)!

### 3 Recurrent Neural Networks: Language Modeling (25 points)

In this section, you'll analyze a recurrent neural network (RNN) used for language modeling.

Language modeling is a central task in NLP, and language models can be found at the heart of speech recognition, machine translation, and many other systems. Given a sequence of words (represented as one-hot vectors)  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$ , a language model predicts the next word  $\mathbf{x}^{(t+1)}$  by modeling:

$$P(\mathbf{x}^{(t+1)} = \mathbf{w}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where  $\mathbf{w}_j$  is a word in the vocabulary.

Your job is to compute the gradients of a recurrent neural network language model, which uses a hidden layer to represent the “history”  $\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}$  of previous words. Formally, the model<sup>6</sup> is:

$$\begin{aligned} \mathbf{e}^{(t)} &= \mathbf{E}\mathbf{x}^{(t)} \\ \mathbf{h}^{(t)} &= \text{sigmoid}(\mathbf{W}_h\mathbf{h}^{(t-1)} + \mathbf{W}_e\mathbf{e}^{(t)} + \mathbf{b}_1) \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2) \\ \bar{P}(\mathbf{x}^{(t+1)} = \mathbf{w}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) &= \hat{y}_j^{(t)} \end{aligned}$$

where  $\mathbf{h}^{(0)} = \mathbf{h}_0 \in \mathbb{R}^{D_h}$  is some initialization vector for the hidden layer and  $\mathbf{E}\mathbf{x}^{(t)}$  is the product of  $\mathbf{E}$  with the one-hot vector  $\mathbf{x}^{(t)}$  representing the current word. Essentially, the first equation embeds the current word, the second equation produces a new hidden state given the current word's embedding and previous hidden state, and the last equation predicts what the next word will be. The parameters are:

$$\mathbf{E} \in \mathbb{R}^{d \times |V|} \quad \mathbf{W}_h \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{W}_e \in \mathbb{R}^{D_h \times d} \quad \mathbf{b}_1 \in \mathbb{R}^{D_h} \quad \mathbf{U} \in \mathbb{R}^{|V| \times D_h} \quad \mathbf{b}_2 \in \mathbb{R}^{|V|} \quad (1)$$

where  $\mathbf{E}$  is the embedding matrix,  $\mathbf{W}_e$  the input word representation matrix,  $\mathbf{W}_h$  the hidden state transformation matrix, and  $\mathbf{U}$  is the output word representation matrix.  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are biases.  $d$  is the embedding size,  $|V|$  is the vocabulary size, and  $D_h$  is the hidden layer size.

The output vector  $\hat{\mathbf{y}}^{(t)} \in \mathbb{R}^{|V|}$  is a probability distribution over the vocabulary. The model is trained by minimizing the (un-regularized) cross-entropy loss:

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$$

<sup>6</sup>This model is adapted from a paper by Toma Mikolov, et al. from 2010: [http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov\\_interspeech2010\\_IS100722.pdf](http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov_interspeech2010_IS100722.pdf)

where  $\mathbf{y}^{(t)}$  is the one-hot vector corresponding to the target word (which here is equal to  $\mathbf{x}^{(t+1)}$ ). We average the cross-entropy loss across all examples (i.e., words) in a sequence to get the loss for a single sequence.

- (a) (4 points, written) Conventionally, when reporting performance of a language model, we evaluate on *perplexity*, which is defined as:

$$PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\bar{P}(\mathbf{x}_{\text{pred}}^{(t+1)} = \mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$

i.e. the inverse probability of the correct word, according to the model distribution  $\bar{P}$ .

- (i) (2 points) Show how you can derive perplexity from the cross-entropy loss (*Hint: remember that  $\mathbf{y}^{(t)}$  is one-hot!*). **This should be a very short problem - not too perplexing!**
- (ii) (1 point) Now use this relationship between perplexity and cross-entropy to show that minimizing the geometric mean perplexity,  $\left(\prod_{t=1}^T PP(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})\right)^{1/T}$ , is equivalent to minimizing the arithmetic mean cross-entropy loss,  $\frac{1}{T} \sum_{t=1}^T CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})$ , across the training set. (*Hint: for any positive function  $f$ , minimizing  $\log(f)$  is equivalent to minimizing  $f$  itself*)
- (iii) (1 point) Suppose you have a vocabulary of  $|V|$  words and your “language model” works by picking the next word uniformly at random from the vocabulary (mathematically,  $\bar{P}(\mathbf{x}^{(t+1)} = \mathbf{w}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = 1/|V|$  for every word  $\mathbf{w}_j$  in the vocabulary). What would the perplexity  $PP(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})$  for a single word to be? Compute the corresponding cross-entropy loss when  $|V| = 10000$ .
- (b) (7 points, written) Compute the gradients of the loss  $J$  (the cross-entropy loss) with respect to the following model parameters at a single point in time  $t$  (to save a bit of time, you don’t have to compute the gradients with the respect to the biases  $\mathbf{b}_1$  and  $\mathbf{b}_2$ ):

$$\frac{\partial J^{(t)}}{\partial \mathbf{U}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t)}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(t)}$$

where  $\Big|_{(t)}$  denotes the gradient for the appearance of that parameter at time  $t$  (in other words,  $\mathbf{h}^{(t-1)}$  is taken to be fixed, so you don’t need to take it’s derivative when applying the chain rule and backpropagate to earlier timesteps - you’ll do that in part (c)).

Additionally, compute the derivative with respect to the *previous* hidden layer value:

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$$

*Hint: you may want to define the following variables:*

$$\begin{aligned} \mathbf{z}^{(t)} &= \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1 \\ \boldsymbol{\theta}^{(t)} &= \mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2 \end{aligned}$$

and define (and compute the value of) the following error terms:

$$\begin{aligned} \boldsymbol{\delta}_1^{(t)} &= \frac{\partial J}{\partial \boldsymbol{\theta}^{(t)}} \\ \boldsymbol{\delta}_2^{(t)} &= \frac{\partial J}{\partial \mathbf{z}^{(t)}} = \boldsymbol{\delta}_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} \end{aligned}$$

- (c) (7 points, written) Below is a sketch of the network at a single timestep:



Draw the “unrolled” network for 3 timesteps (your picture should show  $h^{(t-3)}$  to  $h^{(t)}$ ). Next, compute the following backpropagation-through-time gradients:

$$\frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t-1)}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_{(t-1)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(t-1)}$$

where  $\Big|_{(t-1)}$  denotes the gradient for the appearance of that parameter at time  $(t-1)$ . Use the error term  $\gamma^{(t-1)} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$  to represent the derivative you computed in the previous part (it should show up in three gradients you compute).

We have to compute multiple gradients with respect to parameters like  $\mathbf{W}_e$  because they are used multiple times in the feed-forward computation. The final gradient for the parameters will be the sum of all their gradients over time (e.g., the gradient  $\frac{\partial J^{(t)}}{\partial \mathbf{W}_e} = \sum_{i=0}^t \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_{(i)}$ ).

- (d) (3 points, written) Given  $\mathbf{h}^{(t-1)}$ , how many operations are required to perform backpropagation for a single timestep (i.e., compute the terms you found in part (b)). Express your answer in big-O notation in terms of the dimensions  $d$ ,  $D_h$  and  $|V|$  (Equation 1). Don’t worry about the gradients you didn’t compute (with respect to  $\mathbf{b}_1$  and  $\mathbf{b}_2$ ), they actually don’t change the answer!

(Hint: You only have to worry about matrix multiplications, the other operations do not change the big-O runtime. Multiplying a vector by an  $n$  by  $m$  matrix takes  $O(nm)$  operations.)

- (e) (3 points, written) Now suppose you have a sequence of  $T$  words. How many operations are required to compute the gradient of the loss with respect to the model parameters across the entire sequence ( $\sum_{t=1}^T J^{(t)}(\theta)$ )? Assume we backpropagate through time all the way to  $t = 0$  for each word.

Hint: Look at the computational graph you drew in part (c). Will we have to pass an error signal (upstream gradient) into  $h^{(0)}$   $T$  times (once for the loss from every word), or can we be more efficient than that? Therefore, how many times will we have to do the single timestep (your answer to (d)) computations?

- (f) (1 point, written) Which term in your big-O expressions is likely the largest? Which layer in the RNN is that term from?

**Bonus** (1 point, written) Given your knowledge of similar models (i.e. word2vec), suggest a way to speed up this part of the computation. Your approach can be an approximation, but you should argue why it’s a good one. The paper “Extensions of recurrent neural network language model” (Mikolov, et al. 2013) may be of interest here.