

---

# CS 224N: Assignment 2

RYAN MCMAHON

THURSDAY 2<sup>ND</sup> MARCH, 2017

---

## Contents

<b>1</b>	<b>Problem 1: Tensorflow Softmax (25 pts)</b>	<b>2</b>
1.1	(a) Softmax in Tensorflow (5 pts)	2
1.2	(b) Cross-Entropy Loss in Tensorflow (5 pts)	2
1.3	(c) Placeholders and Feed Dictionaries (5 pts)	3
1.4	(d) Add Softmax and Add CE Loss (5 pts)	3
1.5	(e) Add Training Optimizer – Gradient Descent (5 pts)	3
<b>2</b>	<b>Problem 2: Neural Transition-Based Dependency Parsing (50 pts)</b>	<b>4</b>
2.1	(a) Parsing Steps by Hand (6 pts)	5
2.2	(b) Number of Steps to Parse (2 pts)	5
2.3	(c) Partial Parse Coding (6 pts)	6
2.4	(d) Minibatch Parsing (6 pts)	6
2.5	(N/A) Description (for actual problem see Subproblem 2.6)	7
2.6	(e) Xavier Initialization (4 pts)	8
2.7	(f) Dropout Regularization (2 pts)	8
2.8	(g) Adam Optimizer (4 pts)	9
2.9	(h) Neural Dependency Parser (20 pts)	10
<b>3</b>	<b>Problem 3: Recurrent Neural Networks: Language Modeling (25 pts)</b>	<b>11</b>
3.1	(a) Perplexity and Cross-Entropy Loss (5 pts)	12
3.2	(b) Gradients <i>w.r.t.</i> the Loss (8 pts)	13
3.3	(c) BPTT (8 pts)	15

## Problem 1: Tensorflow Softmax (25 pts)

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

$$J(\mathbf{W}) = CE(\mathbf{y}, \text{softmax}(\mathbf{x}\mathbf{W})) \quad (1.1)$$

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

### 1.1 (a) Softmax in Tensorflow (5 pts)

Implement the softmax function using TensorFlow in `q1_softmax.py`. Remember that

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (1.2)$$

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

**Answer:**

See code: `~/code/q1_softmax.py`.

### 1.2 (b) Cross-Entropy Loss in Tensorflow (5 pts)

Implement the cross-entropy loss using TensorFlow in `q1_softmax.py`. Remember that

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{N_c} y_i \log(\hat{y}_i) \quad (1.3)$$

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 (rows) of 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`.

**Answer:**

See code: `~/code/q1_softmax.py`.

### 1.3 (c) Placeholders and Feed Dictionaries (5 pts)

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 configurations variables in the code.

**Answer:**

The purpose of placeholder variables and feed dictionaries is, in short, to improve computational efficiency. They allow us to specify the graph structure without supplying the actual data. By doing this, our model can be flexible to different training approaches (e.g., different mini-batch sizes). Specifically, the placeholders stand in for whatever data/values we want to use and the feed dictionaries tell the placeholders what data/values to use.

Also, see code: `~/code/q1_classifier.py`.

### 1.4 (d) Add Softmax and Add CE Loss (5 pts)

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, **not** TensorFlow built-ins.

**Answer:**

See code: `~/code/q1_classifier.py`.

### 1.5 (e) Add Training Optimizer – Gradient Descent (5 pts)

Fill in the implementation for `add_training_op` in `q1_classifier.py`. Explain how TensorFlow's automatic differentiation removes the need for us to define gradients explicitly. Verify that your model is able to fit to synthetic data by running `q1_classifier.py` and making sure that the tests pass.

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

**Answer:**

TensorFlow's automatic differentiation works by representing a given model as a symbolic Directed Acyclic Graph (DAG). Using this DAG, we can break down an arbitrarily complex function into basic mathematical operations (e.g.,  $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $\log$ ,  $\dots$  etc.). Once this is done, TensorFlow is able to numerically calculate the necessary partial derivatives for the gradients using the chain-rule, which alleviates the need for us to define them ourselves.

Also, see code: `~/code/q1_classifier.py`.

## Problem 2: Neural Transition-Based Dependency Parsing (50 pts)

*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 of 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.*

## 2.1 (a) Parsing Steps by Hand (6 pts)

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.

**NOTE:** See assignment for table and tree.

**Answer:**

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
[ROOT, parsed, this]	[sentence, correctly]		SHIFT
[ROOT, parsed, this, sentence]	[correctly]		SHIFT
[ROOT, parsed, sentence]	[correctly]	sentence → this	LEFT-ARC
[ROOT, parsed]	[correctly]	parsed → sentence	RIGHT-ARC
[ROOT, parsed, correctly]	[]		SHIFT
[ROOT, parsed]	[]	parsed → correctly	RIGHT-ARC
[ROOT]	[]	ROOT → parsed	RIGHT-ARC

## 2.2 (b) Number of Steps to Parse (2 pts)

A sentence containing  $n$  words will be parsed in how many steps (in terms of  $n$ )? Briefly explain why.

**Answer:**

Assuming we have a grammatically well-formed sentence and need a first step from ROOT:

$$n + 2 \leq \text{Steps} \leq n^3$$

The lower bound exists because we need: 1) an arc from ROOT to the sentence head, and 2) an arc between a subject and a verb. The  $n$  is derived from all of the necessary shifts. Our upper bound comes from Eisner (1996) – whose algorithm allows sub-sequences to be parsed independently of one another.

## 2.3 (c) Partial Parse Coding (6 pts)

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`.

**Answer:**

See code: `~/code/q2_parser_transitions.py`.

## 2.4 (d) Minibatch Parsing (6 pts)

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 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*  $\neq \emptyset$  **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 parses from *unfinished\_pares* ;

**end**

**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.

**Answer:**

See code: `~/code/q2_parser_transitions.py`.

## 2.5 (N/A) Description (for actual problem see Subproblem 2.6)

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  $0 \leq w_i \leq |V|$  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:

$$\mathbf{x} = [\mathbf{L}_{w_1}, \mathbf{L}_{w_2}, \dots, \mathbf{L}_{w_m}] \in \mathbb{R}^{dm}$$

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

$$\mathbf{h} = \text{ReLU}(\mathbf{x}\mathbf{W} + \mathbf{b}_1) \tag{2.1}$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{h}\mathbf{U} + \mathbf{b}_2) \tag{2.2}$$

(recall that  $\text{ReLU}(z) = \max(z, 0)$ ). We evaluate using cross-entropy loss:

$$J(\theta) = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{N_c} y_i \log \hat{y}_i \tag{2.3}$$

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

---

<sup>1</sup>Chen and Manning (2014)

## 2.6 (e) Xavier Initialization (4 pts)

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

$$\epsilon = \frac{\sqrt{6}}{\sqrt{m+n}} \quad (2.4)$$

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  $W$  and  $U$ .

**Answer:**

See code: `~/code/q2_initialization.py`.

## 2.7 (f) Dropout Regularization (2 pts)

We will regularize our network by applying Dropout.<sup>3</sup> During training this randomly sets units in the hidden layer  $h$  to zero with probability  $p$  and then multiplies  $h$  by a constant  $\gamma$  (dropping different units each minibatch). We can write this as

$$h_{drop} = \lambda d \circ h \quad (2.5)$$

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

$$\mathbb{E}_i[h_{drop}]_i = h_i \quad (2.6)$$

for all  $0 \leq i \leq D_h$ . What must  $\gamma$  equal in terms of  $p$ ? Briefly justify your answer.

**Answer:**

$\gamma$  must be the inverse of  $p$ : e.g.,  $\gamma = \frac{1}{p}$ . We have to multiply by the reciprocal of  $p$  to get back to our expectation since we are only seeing the unit  $p$  times. This is best illustrated with an example. Imagine that  $x = [1, 1, 1, 1, 1, 1]$  and  $p = 0.5$ . If we randomly set half of  $x$  to equal zero a bunch of times,  $\bar{x}$  should equal 0.5, and the same will be true for any individual element in  $x$ . This is obviously not the case as we can see that the “true” mean of  $x = 1$ . We can get back to that “true” value by multiplying our estimates by the inverse of  $p$ :  $p^{-1} = \frac{1}{p} = \frac{2}{1} = 2$ .

<sup>2</sup>This is also referred to as Glorot initialization and was initially described in [Glorot and Bengio \(2010\)](#).

<sup>3</sup>[Srivastava et al. \(2014\)](#)



## 2.8 (g) Adam Optimizer (4 pts)

We will train our model using the Adam optimizer.<sup>4</sup> Recall that standard SGD uses the update rule

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J_{\text{minibatch}}(\theta) \quad (2.7)$$

where  $\theta$  is a vector containing all of the model parameters,  $J$  is the loss function,  $\nabla_{\theta} J_{\text{minibatch}}(\theta)$  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>

(i) First, Adam uses a trick called momentum by keeping track of  $m$ , a rolling average of the gradients:

$$m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta) \quad (2.8)$$

$$\theta \leftarrow \theta - \alpha m \quad (2.9)$$

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  $m$  stops the updates from varying as much. Why might this help with learning?

(ii) Adam also uses adaptive learning rates by keeping track of  $v$ , a rolling average of the magnitudes of the gradients:

$$m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J_{\text{minibatch}}(\theta) \quad (2.10)$$

$$v \leftarrow \beta_2 v + (1 - \beta_2) (\nabla_{\theta} J_{\text{minibatch}}(\theta) \circ \nabla_{\theta} J_{\text{minibatch}}(\theta)) \quad (2.11)$$

$$\theta \leftarrow \theta - \alpha m / \sqrt{v} \quad (2.12)$$

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

**Answer:**

(i): The rolling average,  $m$ , reduces variance in our updates by effectively down-weighting the impact of the current minibatch's gradient. This way, any outlier set of observations should not dramatically affect our estimates, but will still have some impact. Imagine that we have gone through 10 minibatches, all of which have told us that some parameter's weight needs to change by a fair bit for us to reach some minima (i.e., a "steep gradient"). Then we get to the 11<sup>th</sup> minibatch and it says we're going the wrong way with our estimates. Intuitively, it makes more sense that this is an outlier batch and that we are indeed traversing the gradient in the correct direction – based on our previous 10 observations. The  $m$  parameter captures this intuition.

(ii): Since Adam divides  $m$  by  $\sqrt{v}$ , the parameters with the steepest gradients and the smallest variances will have the largest updates. This makes sense because we *know* that we need to update those values. If the variance of the gradients is large, we aren't sure where we need to go, so we don't want to move too far – for fear of going the wrong way.

---

<sup>4</sup>Kingma and Ma (2015)

<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.

## 2.9 (h) Neural Dependency Parser (20 pts)

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 training the model won't take as long.
- This code should run within 1 hour on a CPU.
- When running with `debug=False`, you should be able to get a loss smaller than 0.07 on the train set (by the end of the last epoch) and an Unlabeled Attachment Score larger than 88 on the dev set (with the best-performing model out of all the epochs). For comparison, the model in the original neural dependency parsing paper gets 92.5. 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** (1 point): 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 and why it helps (or hurts!) the model. Some extensions may require tweaking the hyperparameters in `Config` to make them effective.

### Answer:

See code: `~/code/q2_parser_model.py`.

The best UAS I achieved on the development data set was 88.67%. This translated to an unlabeled attachment score of 89.07% on the test set.

### Problem 3: Recurrent Neural Networks: Language Modeling (25 pts)

In this section, you'll compute the gradients of a recurrent neural network (RNN) 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 row vectors)  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$ , a language model predicts the next word  $\mathbf{x}^{(t+1)}$  by modeling:

$$P(\mathbf{x}^{(t+1)} = \mathbf{v}_j | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) \quad (3.1)$$

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

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

$$\mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{L} \quad (3.2)$$

$$\mathbf{h}^{(t)} = \text{sigmoid}(\mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1) \quad (3.3)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2) \quad (3.4)$$

$$\bar{P}(\mathbf{x}^{(t+1)} = \mathbf{v}_j | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = \hat{y}_j^{(t)} \quad (3.5)$$

where  $\mathbf{h}^{(0)} = \mathbf{h}_0 \in \mathbb{R}^{D_h}$  is some initialization vector for the hidden layer and  $\mathbf{x}^{(t)} \mathbf{L}$  is the product of  $\mathbf{L}$  with the one-hot row vector  $\mathbf{x}^{(t)}$  representing the current word. The parameters are:

$$\mathbf{L} \in \mathbb{R}^{|V| \times d} \quad \mathbf{H} \in \mathbb{R}^{D_h \times D_h} \quad (3.6)$$

$$\mathbf{I} \in \mathbb{R}^{d \times D_h} \quad \mathbf{b}_1 \in \mathbb{R}^{D_h} \quad (3.7)$$

$$\mathbf{U} \in \mathbb{R}^{D_h \times |V|} \quad \mathbf{b}_2 \in \mathbb{R}^{|V|} \quad (3.8)$$

where  $\mathbf{L}$  is the embedding matrix,  $\mathbf{I}$  is the input word representation matrix,  $\mathbf{H}$  is the hidden transformation matrix, and  $\mathbf{U}$  is the output word representation matrix.  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are biases.  $d$  is the embedding dimension,  $|V|$  is the vocabulary size, and  $D_h$  is the hidden layer dimension.

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:

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)}) \quad (3.9)$$

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.

<sup>1</sup>Model adapted from Mikolov et al. (2010).

### 3.1 (a) Perplexity and Cross-Entropy Loss (5 pts)

Conventionally, when reporting performance of a language model, we evaluate on perplexity, which is defined as:

$$PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\bar{P}(\mathbf{x}_{pred}^{(t+1)} = \mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}} \quad (3.10)$$

i.e. the inverse probability of the correct word, according to the model distribution  $\bar{P}$ . Show how you can derive perplexity from the cross-entropy loss (Hint: remember that  $\mathbf{y}^{(t)}$  is one-hot!), and thus argue that minimizing the (arithmetic) mean cross-entropy loss will also minimize the (geometric) mean perplexity across the training set. **This should be a very short problem - not too perplexing!**

For a vocabulary of  $|V|$  words, what would you expect perplexity to be if your model predictions were completely random (chosen uniformly from the vocabulary)? Compute the corresponding cross-entropy loss for  $|V| = 10000$ .

**Answer:**

(i) Perplexity can alternatively be defined as  $PP(\mathbf{y}, \hat{\mathbf{y}}) = 2^{J^{(t)}(\theta)}$ . We can work back to this definition from the original perplexity formula, in Equation 3.10, above. I drop the superscript  $(t)$  for simplicity:

$$PP(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{\sum_{j=1}^{|V|} y_j \cdot \hat{y}_j}$$

And because we are operating on a one-hot vector, we can move between summations and products:

$$\begin{aligned} &= \frac{1}{\prod_{j=1}^{|V|} \hat{y}_j^{y_j}} \\ &= \frac{1}{\prod_{j=1}^{|V|} 2^{\log(\hat{y}_j^{y_j})}} \\ &= \frac{1}{\prod_{j=1}^{|V|} 2^{y_j \times \log(\hat{y}_j)}} \\ &= 2^{-\sum_{j=1}^{|V|} y_j \cdot \log(\hat{y}_j)} \end{aligned}$$

The reason that minimizing the arithmetic mean cross-entropy loss also minimizes the geometric mean perplexity is that the perplexity can be defined as a monotonic transformation the cross-entropy. In other words,  $J^{(t)}(\theta)$ , and by definition  $2^{J^{(t)}(\theta)}$ , only decreases as  $\hat{y}_{target}^{(t)}$  increases. Thus, minimizing one will also minimize the other. This equivalence is shown below (for  $N$  examples):

$$2^{\frac{1}{N} \sum_{t=1}^{n-1} J^{(t)}(\theta)} = \left( \prod_{t=1}^{n-1} PP(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{\frac{1}{N}}$$

(ii) Assuming  $|V| = 10,000$ , with predictions drawn randomly from a uniform distribution prediction:

$$\begin{aligned}
 E(\hat{y}_j^{(t)}) &= \frac{1}{|V|} = \frac{1}{10,000}, \quad \forall j \in |V| \\
 CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) &= - \sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)}) \\
 &= -\log(E(\hat{y}_{target}^{(t)})) = -\log\left(\frac{1}{10,000}\right) \\
 &\approx 9.213
 \end{aligned}$$

### 3.2 (b) Gradients *w.r.t.* the Loss (8 pts)

Compute the gradients of the loss  $J$  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 respect to  $\mathbf{U}$  and  $\mathbf{b}_1$ ):

$$\frac{\partial J^{(t)}}{\partial \mathbf{b}_2} \quad \frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{I}} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{H}} \Big|_{(t)} \quad (3.11)$$

where  $\mathbf{L}_{x^{(t)}}$  is the row of  $\mathbf{L}$  corresponding to the current word  $x^{(t)}$ , and  $\Big|_{(t)}$  denotes the gradient for the appearance of that parameter at time  $t$  (equivalently,  $\mathbf{h}^{(t-1)}$  is taken to be fixed, and you need not backpropagate to earlier timesteps just yet – you'll do that in part (c)). Additionally, compute the derivative with respect to the previous hidden layer value:

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} \quad (3.12)$$

**Answer:**

Let  $\mathbf{f}_1^{(t)} = \mathbf{h}^{(t-1)}\mathbf{H} + e^{(t)}\mathbf{I} + \mathbf{b}_1$  and  $\mathbf{f}_2^{(t)} = \mathbf{h}^{(t)}\mathbf{U} + \mathbf{b}_2$

And then define the deltas as:

$$\begin{aligned}
 \delta_2^{(t)} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_2^{(t)}} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \\
 \delta_1^{(t)} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_1^{(t)}} = \delta_2^{(t)} \cdot \frac{\partial \mathbf{f}_2^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{f}_1^{(t)}} \\
 &= \delta_2^{(t)} \cdot \mathbf{U}^T \circ \mathbf{h}^{(t)} \circ (1 - \mathbf{h}^{(t)})
 \end{aligned}$$

where the last portion of  $\delta_1^{(t)}$  concerning the item-wise product comes from the earlier derived result, that  $\frac{d}{dx}\sigma(x) = \sigma(x) \times (1 - \sigma(x))$ .

We can then use these  $\delta$  terms to compute the gradients defined in the problem.<sup>2</sup> First, the gradient with respect to  $\mathbf{b}_2$ :

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{b}_2} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_2^{(t)}} \frac{\partial \mathbf{f}_2^{(t)}}{\partial \mathbf{b}_2} \\ &= \delta_2^{(t)} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}\end{aligned}$$

Then with respect to  $\mathbf{L}_{x^{(t)}}$ :

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_1^{(t)}} \frac{\partial \mathbf{f}_1^{(t)}}{\partial \mathbf{e}^{(t)}} \frac{\partial \mathbf{e}^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} \\ &= \delta_1^{(t)} \mathbf{I}^T\end{aligned}$$

And the with respect to  $\mathbf{I}$ :

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{I}} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_1^{(t)}} \frac{\partial \mathbf{f}_1^{(t)}}{\partial \mathbf{I}} \\ &= (\mathbf{e}^{(t)})^T \delta_1^{(t)}, \text{ note at } \ddagger\end{aligned}$$

$\ddagger$  : Originally had this as  $\left((\delta_1^{(t)})^T \mathbf{e}^{(t)}\right)^T$ , but unnecessarily transposes two elements. Solution shown above is much cleaner.

Now with respect to  $\mathbf{H}$ :

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{H}} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_1^{(t)}} \frac{\partial \mathbf{f}_1^{(t)}}{\partial \mathbf{H}} \\ &= (\mathbf{h}^{(t-1)})^T \delta_1^{(t)}\end{aligned}$$

And finally with respect to  $\mathbf{h}^{(t-1)}$ :

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{h}^{(t-1)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{f}_1^{(t)}} \frac{\partial \mathbf{f}_1^{(t)}}{\partial \mathbf{h}^{(t-1)}} \\ &= \delta_1^{(t)} \mathbf{H}^T\end{aligned}$$

---

<sup>2</sup>I should note that I did have some assistance from the solutions document here – mainly with the definition of  $\delta_1^{(t)}$  above: when the cross product was applicable versus item-wise products.

### 3.3 (c) BPTT (8 pts)

Draw the unrolled network for 3 timesteps, and compute the backpropagation-through-time gradients:

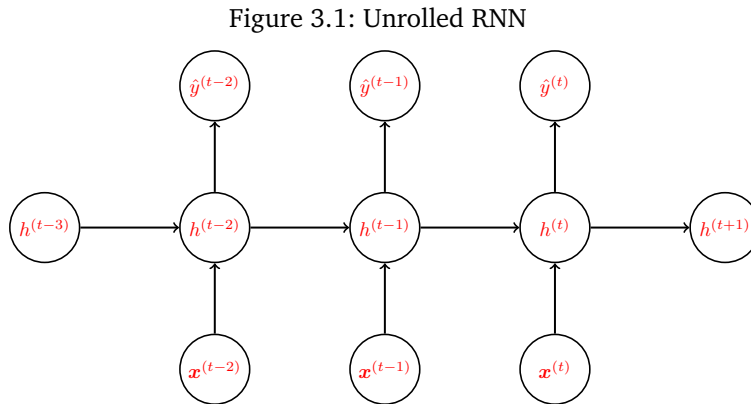
$$\frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t-1)}}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{I}} \Big|_{(t-1)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{H}} \Big|_{(t-1)} \quad (3.13)$$

where  $\Big|_{(t-1)}$  denotes the gradient for the appearance of that parameter at time  $(t-1)$ . Because parameters are used multiple times in feed-forward computation, we need to compute the gradient for each time they appear.

You should use the backpropagation rules from Lecture 5 to express these derivatives in terms of error term  $\delta^{(t-1)} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$  computed in the previous part. (Doing so will allow for re-use of expressions for  $t-2$ ,  $t-3$ , and so on.

Note that the true gradient with respect to a training example requires us to run backpropagation all the way back to  $t=0$ . In practice, however, we generally truncate this and only backpropagate for a fixed number  $\tau \approx 5-10$  timesteps.

**Answer:**



Now for the gradient calculations, let  $\sigma'(\mathbf{f}_1^{(t-1)}) = \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{f}_1^{(t-1)}} = \text{diag}(\mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}))$