

CS 224n: Assignment #2

Due date: 2/9 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 the course website.

Note: In this assignment, the inputs to neural network layers will be row vectors because this is standard practice for TensorFlow (some built-in TensorFlow functions assume the inputs are row vectors). This means the weight matrix of a hidden layer will right-multiply instead of left-multiply its input (i.e., $\mathbf{x}\mathbf{W} + \mathbf{b}$ instead of $\mathbf{W}\mathbf{x} + \mathbf{b}$).

1 Tensorflow Softmax (25 points)

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

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

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.

- (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 (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`.

- (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.

Solution: Placeholder variables and feed dictionaries make it possible to feed data (such as training examples for a neural network) into the computational graph.

- (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, **not** TensorFlow built-ins.
- (e) (5 points, coding/written) 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 `python q1_classifier.py` and making sure that the tests pass.

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

Solution: TensorFlow's automatic differentiation means we only need to define the forward pass of our model; the backwards pass is done automatically, allowing us to optimize the model without explicitly defining gradients.

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:

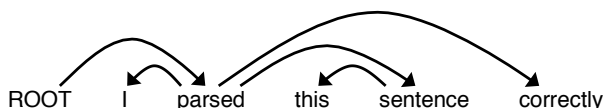
- A *stack* of words that are currently being processed.
- A *buffer* of words yet to be processed.
- A list of *dependencies* predicted by the parser.

Initially, the stack only contains `ROOT`, the dependencies list is empty, and the buffer contains all words of the sentence in order. At each step, the parser applies a *transition* to the partial parse until its buffer is empty and only `ROOT` is on the stack. 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 |

Solution:

| stack | buffer | new dependency | transition |
|--------------------------------|-----------------------|------------------|------------|
| [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 |

- (b) (2 points, written) A sentence containing n words will be parsed in how many steps (in terms of n)? Briefly explain why.

Solution: Each word of the sentence must be shifted onto the stack and then reduced away, so a sentence containing n words will be parsed in $2n$ steps.

- (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_parses` as a list of partial parses, one for each sentence in `sentences`

Initialize `unfinished_parses` as a shallow copy of `partial_parses`

while `unfinished_parses` is not empty **do**

 Take the first `batch_size` parses in `unfinished_parses` 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_parses`

end while

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

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*¹. 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:

$$\mathbf{x} = [\mathbf{L}_{w_0}, \mathbf{L}_{w_1}, \dots, \mathbf{L}_{w_m}] \in \mathbb{R}^{dm}$$

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:

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

(recall that $\text{ReLU}(z) = \max(z, 0)$). We evaluate using 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 initialization².

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} .

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

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

- (f) (2 points, written) We will regularize our network by applying Dropout³. 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 γ (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 $(1 - p_{drop})$. γ 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 γ equal in terms of p_{drop} ? Briefly justify your answer.

Solution:

$$\mathbb{E}_{p_{drop}}[\mathbf{h}_{drop}]_i = \mathbb{E}_{p_{drop}}[\gamma d_i h_i] = p_{drop}(0) + (1 - p_{drop})\gamma h_i = (1 - p_{drop})\gamma h_i = h_i$$

So γ must equal $1/(1 - p_{drop})$

- (g) (4 points, written) We will train our model using the Adam⁴ 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 α is the learning rate. Adam uses a more sophisticated update rule with two additional steps⁵.

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

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

where β_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?

Solution: Each update will be mostly the same as the previous one (only $1 - \beta_1$ of \mathbf{m} changes each step), so the updates won't vary as much. One way of thinking about this is that it will stop the model parameters from "bouncing around as much" when moving towards a local optimum. Another way is that doing the rolling average is a bit like computing the gradient over a larger minibatch, so each update will be closer to the true gradient over the whole dataset (i.e., lower variance means each gradient estimate is closer to the mean).

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

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

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

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

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

where \circ and $/$ denote elementwise multiplication and division (so $\mathbf{z} \circ \mathbf{z}$ is elementwise squaring) and β_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?

Solution: The parameters with the smallest gradients (on average) will get the larger updates. This means parameters that are at a place where the loss with respect to them is pretty flat will get larger updates, helping them move off plateaus.

- (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.
- You should be able to get a loss smaller than 0.07 on the train set and an Unlabeled Attachment Score larger than 88 on the dev set. 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.

3 Recurrent Neural Networks: Language Modeling (25 points)

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:

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

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 model⁶ is, for $t = 1, \dots, n - 1$:

$$\begin{aligned} \mathbf{e}^{(t)} &= \mathbf{x}^{(t)} \mathbf{L} \\ \mathbf{h}^{(t)} &= \text{sigmoid} \left(\mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1 \right) \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax} \left(\mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2 \right) \\ \bar{P}(\mathbf{x}^{(t+1)} = \mathbf{v}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) &= \hat{y}_j^{(t)} \end{aligned}$$

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

$$\mathbf{L} \in \mathbb{R}^{|V| \times d} \quad \mathbf{H} \in \mathbb{R}^{D_h \times D_h} \quad \mathbf{I} \in \mathbb{R}^{d \times D_h} \quad \mathbf{b}_1 \in \mathbb{R}^{D_h} \quad \mathbf{U} \in \mathbb{R}^{D_h \times |V|} \quad \mathbf{b}_2 \in \mathbb{R}^{|V|} \quad (1)$$

where \mathbf{L} is the embedding matrix, \mathbf{I} the input word representation matrix, \mathbf{H} 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:

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

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

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

i.e. the inverse probability of the correct word, according to the model distribution \bar{P} . 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 $|V| = 10000$.

Solution: Since $\mathbf{y}^{(t)}$ is one-hot, suppose without loss of generality that $y_i^{(t)}$ is the only nonzero element of $\mathbf{y}^{(t)}$. Then

$$\begin{aligned} CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) &= -\log \hat{y}_i^{(t)} = \log \frac{1}{\hat{y}_i^{(t)}} \\ \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) &= \frac{1}{\hat{y}_i^{(t)}} \end{aligned}$$

⁶This 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

It follows that

$$\text{CE}(y^{(t)}, \hat{y}^{(t)}) = \log \text{PP}^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

Thus minimizing the arithmetic mean of the cross-entropy is identical to minimizing the geometric mean of the perplexity. If the model predictions are completely random, $E[\hat{y}_i^{(t)}] = \frac{1}{|V|}$. Given that $y^{(t)}$ is one-hot, the expected value of the perplexity is $1/(1/|V|) = |V|$. Since the cross-entropy is the logarithm of perplexity, the expected cross-entropy if $|V| = 10000$ is $\log 10000 \approx 9.21$.

- (b) (8 points, written) 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 the respect to \mathbf{U} and \mathbf{b}_1):

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

where $\mathbf{L}_{x^{(t)}}$ is the row of \mathbf{L} corresponding to the current word $\mathbf{x}^{(t)}$, 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:

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

Solution:

For convenience, let

$$\begin{aligned} \mathbf{v}^{(t)} &= \mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1 \\ \boldsymbol{\theta}^{(t)} &= \mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2 \end{aligned}$$

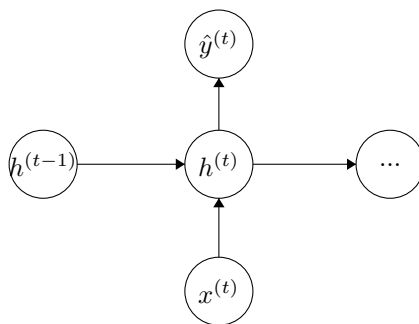
Recall that $\frac{d}{dz} \text{sigmoid}(z) = \text{sigmoid}(z)(1 - \text{sigmoid}(z))$. We can define the error signals

$$\begin{aligned} \delta_1^{(t)} &= \frac{\partial J}{\partial \boldsymbol{\theta}^{(t)}} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \\ \delta_2^{(t)} &= \frac{\partial J}{\partial \mathbf{v}^{(t)}} = \delta_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{v}^{(t)}} = \delta_1^{(t)} \mathbf{U}^T \circ \mathbf{h}^{(t)} \circ (1 - \mathbf{h}^{(t)}) \end{aligned}$$

So

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{b}_2} &= \frac{\partial J}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{b}_2} = \delta_1^{(t)} \\ \frac{\partial J}{\partial \mathbf{L}_{x^{(t)}}} &= \frac{\partial J}{\partial \mathbf{v}^{(t)}} \frac{\partial \mathbf{v}^{(t)}}{\partial \mathbf{e}^{(t)}} \frac{\partial \mathbf{e}^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} = \delta_2^{(t)} \mathbf{I}^T \\ \frac{\partial J}{\partial \mathbf{I}} &= \frac{\partial J}{\partial \mathbf{v}^{(t)}} \frac{\partial \mathbf{v}^{(t)}}{\partial \mathbf{I}} = (\mathbf{e}^{(t)})^T \delta_2^{(t)} \\ \frac{\partial J}{\partial \mathbf{H}} &= \frac{\partial J}{\partial \mathbf{v}^{(t)}} \frac{\partial \mathbf{v}^{(t)}}{\partial \mathbf{H}} = (\mathbf{h}^{(t-1)})^T \delta_2^{(t)} \\ \frac{\partial J}{\partial \mathbf{h}^{(t-1)}} &= \frac{\partial J}{\partial \mathbf{v}^{(t)}} \frac{\partial \mathbf{v}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \delta_2^{(t)} \mathbf{H}^T \end{aligned}$$

(c) (8 points, written) Below is a sketch of the network at a single timestep:



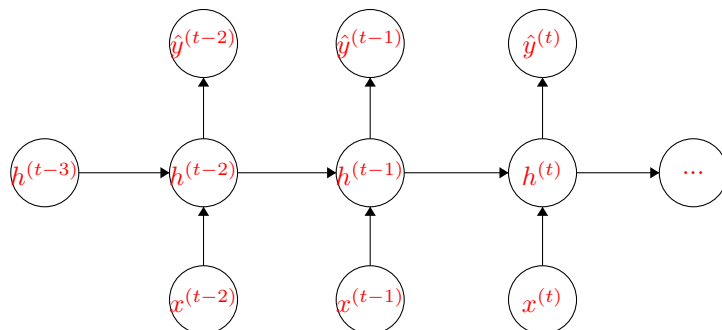
Draw the “unrolled” network for 3 timesteps, and compute the backpropagation-through-time gradients:

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

where $\left|_{(t-1)}\right.$ 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 5⁷ 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.



Solution:

$$\text{Let } \sigma'(\mathbf{v}^{(t-1)}) = \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{v}^{(t-1)}} = \text{diag}(\mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}))$$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{L}_{\mathbf{x}^{(t-1)}}} &= \frac{\partial J}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{v}^{(t-1)}} \frac{\partial \mathbf{v}^{(t-1)}}{\partial \mathbf{L}_{\mathbf{x}^{(t-1)}}} = \boldsymbol{\delta}^{(t-1)} \sigma'(\mathbf{v}^{(t-1)}) \mathbf{I}^T \\ \frac{\partial J}{\partial \mathbf{I}} &= \frac{\partial J}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{v}^{(t-1)}} \frac{\partial \mathbf{v}^{(t-1)}}{\partial \mathbf{I}} = (\mathbf{e}^{(t-1)})^T \boldsymbol{\delta}^{(t-1)} \sigma'(\mathbf{v}^{(t-1)}) \\ \frac{\partial J}{\partial \mathbf{H}} &= \frac{\partial J}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{v}^{(t-1)}} \frac{\partial \mathbf{v}^{(t-1)}}{\partial \mathbf{H}} = (\mathbf{h}^{(t-2)})^T \boldsymbol{\delta}^{(t-1)} \sigma'(\mathbf{v}^{(t-1)}) \end{aligned}$$

(d) (4 points, written) Given $\mathbf{h}^{(t-1)}$, how many operations are required to perform one step of forward propagation to compute $J^{(t)}(\theta)$? How about backpropagation for a single step in time? For τ steps in

⁷<https://web.stanford.edu/class/cs224n/lectures/cs224n-2017-lecture5.pdf>

time? Express your answer in big-O notation in terms of the dimensions d , D_h and $|V|$ (Equation 1). What is the slow step?

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.

Solution: **Forward prop:** $O(|V|D_h + dD_h + D_h^2)$

Backward prop: $O(\tau(|V|D_h + dD_h + D_h^2))$

This question could be interpreted as as back-propping from a single word instead of all words in the τ timesteps. In this case the complexity is $O(|V|D_h + \tau(dD_h + D_h^2))$. However, it's worth mentioning that in practice doing bptt from each word one at a time is inefficient (it would take $O(\tau|V|D_h + \tau^2(dD_h + D_h^2))$ operations).

Slow part is the $O(|V|D_h)$ term from the matrix multiply when computing a probability distribution over next words (assuming $|V| \gg D_h$).