

CS224N: NATURAL LANGUAGE PROCESSING WITH DEEP LEARNING
ASSIGNMENT #2

ANTHONY HO

1. (a) Please see the coding portion of the assignment.
- (b) Please see the coding portion of the assignment.
- (c) The purpose of the placeholder variables is to allocate storage for data/labels before building the computation graph. The feed dictionaries allows us to inject data/labels into the placeholders in a computation graph. Please see the coding portion of the assignment for implementation.
- (d) Please see the coding portion of the assignment.
- (e) When the model's `train_op` is called, (1) it creates a gradient descent optimizer; (2) it calls `add_loss_op` to compute the cross entropy loss based on the data, labels, and current values of the variables `W` and `b`; (3) it computes the gradients w.r.t the loss via automatic differentiation; (4) and at the end it updates the values of the variables `W` and `b` in the direction of the gradient and in proportion to the learning rate as defined in `Config`.
Please see the coding portion of the assignment for implementation.

2. (a) The sequence of transitions are:

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

- (b) A sentence containing n words will be parsed in $2n$ steps, since each word must be first shifted from the buffer into the stack and then removed from the stack as a dependent of another item.
- (c) Please see the coding portion of the assignment.
- (d) Please see the coding portion of the assignment.
- (e) Please see the coding portion of the assignment.
- (f) For the following equation to be true:

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

γ must fulfill the following criteria:

$$\begin{aligned} \mathbb{E}_{p_{drop}}[\mathbf{h}_{drop}]_i &= h_i \\ \implies \gamma(1 - p_{drop})h_i &= h_i \\ \implies \gamma &= \frac{1}{1 - p_{drop}} \end{aligned}$$

- (g) (i) By using \mathbf{m} and a β_1 of 0.9, the new $\boldsymbol{\theta}$ would only be updated slightly towards the new direction and would be largely the same as the previous $\boldsymbol{\theta}$. It helps the updates in $\boldsymbol{\theta}$ to maintain a relatively steady trajectory and prevents the updates from “diffusing around” too much, and thus helps speeding up reaching the local optimum.
- (ii) Since Adam divides the updates by $\sqrt{\mathbf{v}}$, the model parameters that have smaller magnitudes will get larger updates. This might help with combating the “saturated neurons” problem by giving a “boost” to the updates of the parameters that are “saturated” to get out of the “plateaus”.
- (h) Please see the coding portion of the assignment for implementation. The best UAS achieved on the dev set is 88.66 and the UAS achieved on the test is 89.17.
3. (a) (i) Let’s denote k as the index for the target word. Since $\mathbf{y}^{(t)}$ is a one-hot vector:

$$\text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\sum_{j=1}^{|V|} y_i^{(t)} \cdot \hat{y}_i^{(t)}} = \frac{1}{\hat{y}_k^{(t)}} \quad (1)$$

and:

$$CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{j=1}^{|V|} y_i^{(t)} \log \hat{y}_i^{(t)} = - \log \hat{y}_k^{(t)} \quad (2)$$

Therefore, combining equation (1) and (2):

$$CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \log \hat{y}_k^{(t)} = - \log \frac{1}{\text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})} = \log \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \quad (3)$$

- (ii) We can rewrite the log of geometric mean perplexity using equation (3):

$$\begin{aligned} \log \left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{1/T} &= \frac{1}{T} \log \left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right) \\ &= \frac{1}{T} \sum_{t=1}^T \log \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \\ &= \frac{1}{T} \sum_{t=1}^T CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \end{aligned}$$

Since $\left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{1/T}$ is a positive function, minimizing $\log \left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{1/T}$ is equivalent to minimizing $\left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{1/T}$ itself. Therefore, minimizing the geometric mean perplexity $\left(\prod_{t=1}^T \text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right)^{1/T}$ is equivalent to minimizing the arithmetic mean cross-entropy loss $\frac{1}{T} \sum_{t=1}^T CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})$.

- (iii) If $\bar{P}(\mathbf{x}_{\text{pred}}^{(t+1)} = \mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = 1/|V|$, it means $\text{PP}^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = 1/(1/|V|) = |V|$. When $|V| = 10000$, the corresponding cross-entropy loss is $\log(10000) = 9.21$.

- (b) Let’s denote:

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

We can define and compute the values of the following error terms:

$$\begin{aligned} \boldsymbol{\sigma}_1^{(t)} &= \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} = \frac{\partial CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})}{\partial \boldsymbol{\theta}^{(t)}} = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \in \mathbb{R}^{|V| \times 1} \\ \boldsymbol{\sigma}_2^{(t)} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} = \boldsymbol{\sigma}_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} = \mathbf{U}^\top \left(\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \right) \circ \mathbf{h}^{(t)} \circ (1 - \mathbf{h}^{(t)}) \in \mathbb{R}^{D_h \times 1} \end{aligned}$$

Therefore,

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{U}} &= \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{U}} = \boldsymbol{\sigma}_1^{(t)} \left(\mathbf{h}^{(t)} \right)^\top \in \mathbb{R}^{|V| \times D_h} \\ \frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{e}^{(t)}} = \mathbf{W}_e^\top \boldsymbol{\sigma}_2^{(t)} \in \mathbb{R}^{d \times 1} \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} &= \left. \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} = \boldsymbol{\sigma}_2^{(t)} \left(\mathbf{e}^{(t)} \right)^\top \in \mathbb{R}^{D_h \times d} \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} &= \left. \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} = \boldsymbol{\sigma}_2^{(t)} \left(\mathbf{h}^{(t-1)} \right)^\top \in \mathbb{R}^{D_h \times D_h} \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \mathbf{W}_h^\top \boldsymbol{\sigma}_2^{(t)} \in \mathbb{R}^{D_h \times 1}\end{aligned}$$

(c) The “unrolled” network for 3 timesteps is shown in figure 1.

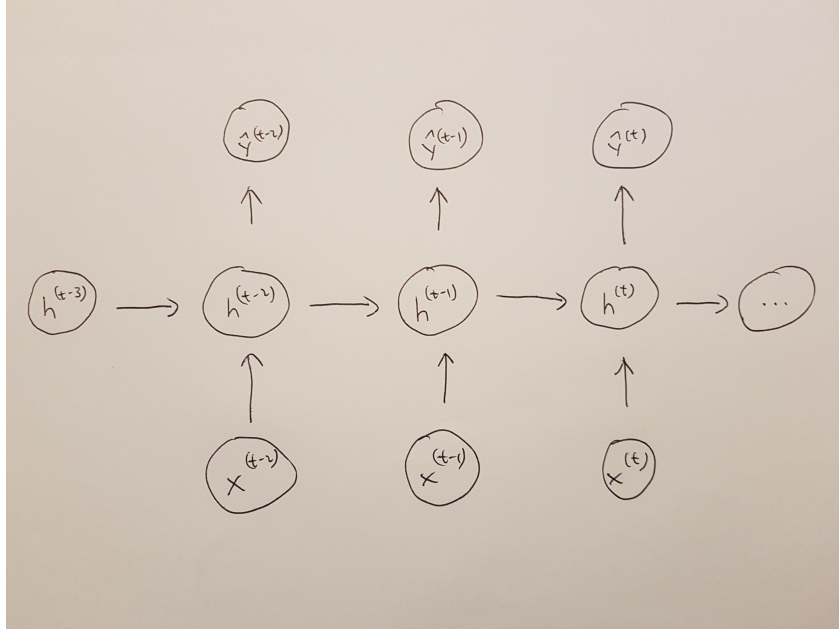


FIGURE 1. The “unrolled” network for 3 timesteps.

Let’s compute the gradient of $J^{(t)}$ w.r.t. $\mathbf{z}^{(t-1)}$:

$$\frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t-1)}} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{z}^{(t-1)}} = \boldsymbol{\gamma}^{(t-1)} \circ \mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}) \in \mathbb{R}^{D_h \times 1}$$

Therefore,

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t-1)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{e}^{(t-1)}} = \mathbf{W}_e^\top \left(\boldsymbol{\gamma}^{(t-1)} \circ \mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}) \right) \in \mathbb{R}^{d \times 1} \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \right|_{(t-1)} &= \left. \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{W}_e} \right|_{(t-1)} = \left(\boldsymbol{\gamma}^{(t-1)} \circ \mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}) \right) \left(\mathbf{e}^{(t-1)} \right)^\top \in \mathbb{R}^{D_h \times d} \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right|_{(t-1)} &= \left. \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{W}_h} \right|_{(t-1)} = \left(\boldsymbol{\gamma}^{(t-1)} \circ \mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}) \right) \left(\mathbf{h}^{(t-2)} \right)^\top \in \mathbb{R}^{D_h \times D_h}\end{aligned}$$

(d) To perform backpropagation for a single timestep, we will have to compute $\frac{\partial J^{(t)}}{\partial \mathbf{U}}$, $\left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)}$, $\left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)}$, $\frac{\partial J^{(t)}}{\partial \mathbf{b}_1}$, and $\frac{\partial J^{(t)}}{\partial \mathbf{b}_2}$. Since $\frac{\partial J^{(t)}}{\partial \mathbf{b}_2} = \boldsymbol{\sigma}_2^{(t)}$ and $\frac{\partial J^{(t)}}{\partial \mathbf{b}_1} = \boldsymbol{\sigma}_1^{(t)}$, their computations are implicit to the computations of $\frac{\partial J^{(t)}}{\partial \mathbf{U}}$, $\left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)}$, and $\left. \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)}$ and thus won’t contribute to the overall time complexity. The rest of the terms that involve matrix multiplications are as follows:

Term	Matrix multiplication	Time complexity
$\sigma_2^{(t)}$	$\mathbf{U}^\top (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)})$	$\mathcal{O}(D_h V)$
$\frac{\partial J^{(t)}}{\partial \mathbf{U}}$	$\sigma_1^{(t)} (\mathbf{h}^{(t)})^\top$	$\mathcal{O}(D_h V)$
$\frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big _{(t)}$	$\sigma_2^{(t)} (\mathbf{e}^{(t)})^\top$	$\mathcal{O}(D_h d)$
$\frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big _{(t)}$	$\sigma_2^{(t)} (\mathbf{h}^{(t-1)})^\top$	$\mathcal{O}(D_h^2)$

Therefore, given $\mathbf{h}^{(t-1)}$, the time complexity of performing backpropagation for a single timestep is:

$$\mathcal{O}(D_h|V| + D_h d + D_h^2)$$

- (e) Since we can store and reuse $\gamma^{(t-i)}$ at each timestep, we will just have to repeat the single timestep in part (d) by T times. Therefore, the time complexity of computing the gradient of the loss w.r.t. the model parameters across the entire sequence is:

$$\mathcal{O}(T(D_h|V| + D_h d + D_h^2))$$

- (f) The term $D_h|V|$ is likely the largest, since presumably the size of the vocabulary $|V|$ is much larger than hidden layer size D_h or the embedding size d . The term corresponds to the RNN layer that computes $\hat{\mathbf{y}}^{(t)}$, i.e. the probabilities of each of the word in the vocabulary being the next word.