

## Assignment #2

ZIAUL  
HASAN  
HASNMI

### ① Tensorflow Softmax

(a) Please check code

(b) Please check code

(c) Placeholder variables allow us to assign data at a later stage hence create operations to build the computational graph without requiring the data beforehand. Feed dictionaries make it possible to feed data such as training data into the computational graph.

Please check code for implementation

(d) Please check code for implementation

(e) Please check code for implementation part

When `train_op` is called, Tensorflow automatically computes the gradients of all the variables (weights & biases) in the computational graph by minimizing the loss via gradient descent algorithm.

## ② Neural Transition - Based Dependency Parsing

(a)

stack	buffer	new dependency	transition
[ROOT]	[I, parsed, this, sentence correctly]		Initial Config
[ROOT, I]	[parsed, this, sentence, correctly]		Shift
[ROOT, I, Parsed]	[this, sentence, correctly]		Shift
[ROOT, Parsed]	[this, sentence, correctly]	parsed → I	Left-ARG
[ROOT, Parsed, this]	[sentence, correctly]		Shift
[ROOT, Parsed, this, sentence]	[correctly]		Shift
[ROOT, Parsed, sentence]	[correctly]	Sentence → this	Left-ARG
[ROOT, Parsed]	[correctly]	Parsed → sentence	Right-ARC
[ROOT, Parsed, correctly]	∅		Shift
[ROOT, Parsed]	∅	Parsed - correctly	Right-ARC
[ROOT]	∅	ROOT - Parsed	Left-ARC

(b) Since each word is processed twice; once to move to the stack and once to create dependency, so for a sentence containing  $n$  words we need " $2n$ " steps.

(c) Please check code implementation

(d) Please check code implementation

(e) Please check code for implementation

(f)  $E_{\text{drop}}[h_{\text{drop}}]_i$ :

$$= \gamma [\hat{p}_{\text{drop}}(0) + (1 - \hat{p}_{\text{drop}})(h_i)]$$

$$= \gamma (1 - \hat{p}_{\text{drop}}) h_i$$

If  $E_{\text{drop}}[h_{\text{drop}}]_i = h_i$

$$\Rightarrow \gamma = \frac{1}{1 - \hat{p}_{\text{drop}}}$$

(g) (i) Each "m" will approximately be the same as previous one because  $(1 - \beta_1)$  is small. Therefore updates won't vary that much.

Since minibatch gradients are random approximations of the true gradient, momentum calculates a rolling average of gradients and stops the oscillations. This helps SGD to accelerate in the right direction. Essentially stops parameters from bouncing around too much.

(ii) The model parameters with smaller average gradients will get larger updates. This should help to move such parameters for which loss doesn't change much, effectively moving them from flat portions (plateaus).

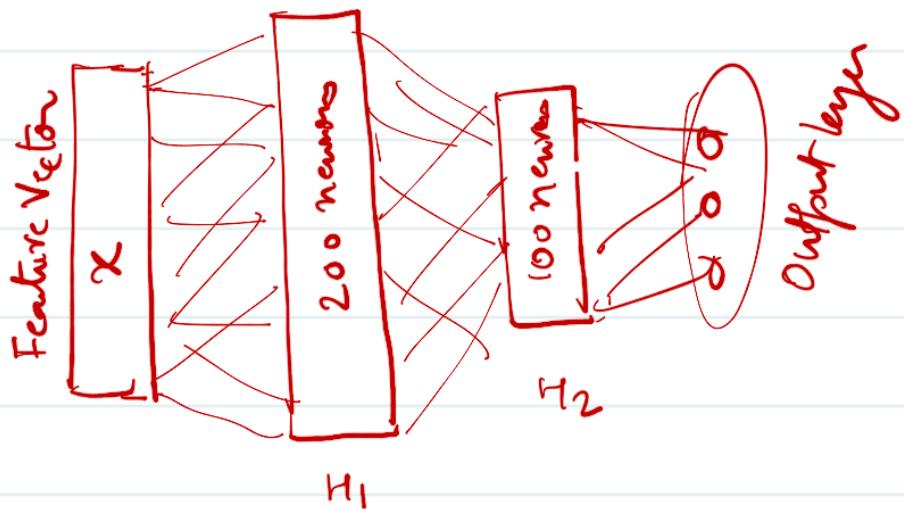
(h) **\*** Please check the implementation of my parser  
in "q2-parser-model.py"

**\*** dev set UAS = 88.66

test set UAS = 84.24

**\*** Please check q2-test.predicted.pkl file.

(i) I added one additional hidden layer with 100 units (with Dropouts)



I also reduced the dropout probability to 0.2 so keep probability is 0.8 & increased epochs to 12.

Evaluation result:

★ dev set UAS = 88.98

★ test set UAS = 85.39

Adding additional layer helps to compute more complex features that can help learn the dependency structures. Also, I increased epochs to train the parameters even more and reduced dropouts to reduce regularization effect. Perhaps we can play with parameters such as learning rate and hidden layer size to further improve the results.

### ③ Recurrent Neural Networks: Language Modelling

(a) (i) if  $y^{(t)}$  is one hot, assume  $y_i^{(t)} = 1$  and zero everywhere else, then

$$CE(y^{(t)}, \hat{y}^{(t)}) = -\log \hat{y}_i^{(t)} = \log \frac{1}{\hat{y}_i^{(t)}}$$

$$\& PP^{(t)}(y^{(t)}, \hat{y}^{(t)}) = \frac{1}{\hat{y}_i^{(t)}}$$

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

$$\Rightarrow \boxed{PP^{(t)}(y^{(t)}, \hat{y}^{(t)}) = e^{CE(y^{(t)}, \hat{y}^{(t)})}}$$

(ii)

$$\arg \min \frac{1}{T} \sum_{t=1}^T CE(y^{(t)}, \hat{y}^{(t)})$$

$$= \arg \min \frac{1}{T} \sum_{t=1}^T \log p_p^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

$$= \arg \min \sum_{t=1}^T \log \left[ p_p^{(t)}(y^{(t)}, \hat{y}^{(t)}) \right]^{\frac{1}{T}}$$

$$= \arg \min \log \prod_{t=1}^T \left( p_p^{(t)}(y^{(t)}, \hat{y}^{(t)}) \right)^{\frac{1}{T}}$$

$$= \arg \min \prod_{t=1}^T p_p^{(t)}(y^{(t)}, \hat{y}^{(t)})^{\frac{1}{T}}$$

Q.E.D

(ii)

Perplexity  $PP(y^{(t)}, \hat{y}^{(t)})$

$$= \frac{1}{|V|}$$

$$= |V|$$

Perplexity would therefore be  $|V|$  for a single word.

Cross entropy  $CE(y^{(t)}, \hat{y}^{(t)}) = \log(\text{perplexity})$

$$= \log |V|$$

for  $|V| = 10000$

$$CE = 9.21$$

(assuming natural log)

$$(b) e^{(t)} = E x^{(t)}$$

$$h^{(t)} = \text{sigmoid}(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2)$$

As per these equations

$h^{(t)}, y^{(t)}, \hat{y}^{(t)}, e^{(t)}$  are assumed to be column vectors. We assume gradients are row vectors so appropriate transpose is required during gradient updates

Assume  $\boxed{\begin{aligned} Z^{(t)} &= W_h h^{(t-1)} + W_e e^{(t)} + b_1 \\ \theta^{(t)} &= U h^{(t)} + b_2 \end{aligned}}$

$$\delta_1^{(t)} = \frac{\partial J^{(t)}}{\partial \theta^{(t)}} = (\hat{y}^{(t)} - y^{(t)})^T \in \mathbb{R}^{1 \times |V|} \quad [\text{From Assignment #1}]$$

$$\delta_2^{(t)} = \frac{\partial J^{(t)}}{\partial Z^{(t)}} = \delta_1^{(t)} \frac{\partial \theta^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial Z^{(t)}} \quad [\text{since } \sigma'(z) = \sigma(z)(1-\sigma(z))]$$

$$= \delta_1^{(t)} U \circ h^{(t)} \circ (1 - h^{(t)}) \in \mathbb{R}^{1 \times D_h}$$

$$\boxed{\frac{\partial J^{(t)}}{\partial U} = (\tilde{y}^{(t)} - y^{(t)}) (h^{(t)})^T} \in \mathbb{R}^{|V| \times D_h}$$

$|V| \times 1$        $1 \times D_h$

(We derived  
similar derivative  
in Assignment #1)

$$\boxed{\frac{\partial J^{(t)}}{\partial e^{(t)}} = \delta_2^{(t)} \cdot w_e} \in \mathbb{R}^{1 \times d}$$

$1 \times D_h$        $D_h \times d$

(Dimension check)

$$\boxed{\left. \frac{\partial J^{(t)}}{\partial w_e} \right|_t = (\delta_2^{(t)})^T \cdot (e^{(t)})^T} \in \mathbb{R}^{D_h \times d}$$

$(D_h \times 1)$        $(1 \times d)$

(Dimension check)

$$\boxed{\left. \frac{\partial J^{(t)}}{\partial w_h} \right|_t = (\delta_2^{(t)})^T \cdot (h^{(t-1)})^T} \in \mathbb{R}^{D_h \times D_h}$$

$D_h \times 1$        $1 \times D_h$

(Dimension check)

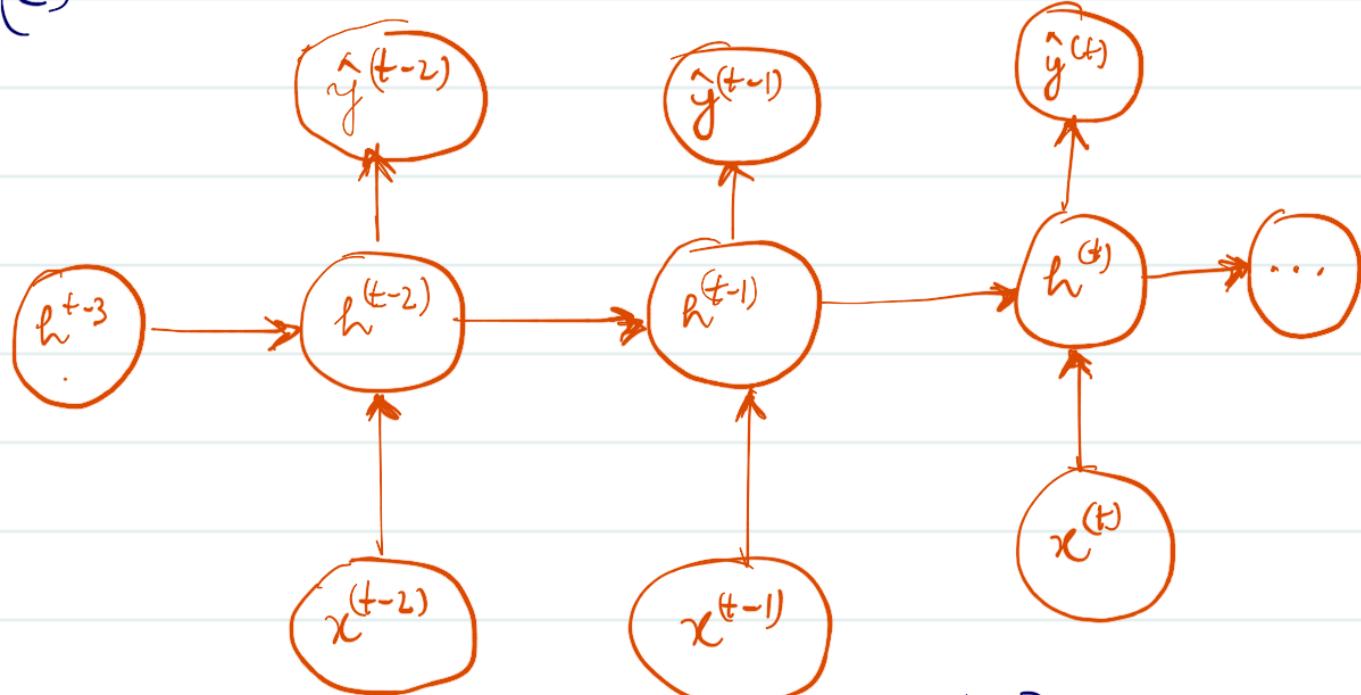
Finally ,

$$\boxed{\frac{\partial J^{(t)}}{\partial h^{(t-1)}} = \delta_2^{(t)} \cdot w_h} \in \mathbb{R}^{1 \times D_w}$$

$1 \times D_h$        $D_h \times D_h$

(Dimension check)

(C)



Let's say  $\gamma^{(t-1)} = \frac{\partial J^{(t)}}{\partial h^{(t-1)}}$   $\in \mathbb{R}^{1 \times D_h}$

$$h^{(t-1)} = \text{Sigmoid}(z^{(t-1)})$$

$$\text{where } z^{(t-1)} = W_h h^{(t-2)} + W_e e^{(t-1)} + b_1$$

$$\frac{\partial h^{(t-1)}}{\partial z^{(t-1)}} = \text{diag}(h^{(t-1)} \circ (1 - h^{(t-1)}))$$

$$\delta_2^{(t-1)} = \frac{\partial J^{(t)}}{\partial e^{(t-1)}} = \gamma^{(t-1)} \circ h^{(t-1)} \circ (1 - h^{(t-1)}) \in \mathbb{R}^{1 \times D_h}$$

$$\therefore \frac{\partial J^{(t)}}{\partial e^{(t-1)}} = \delta_2^{(t-1)} \cdot W_e$$

$$\boxed{\frac{\partial J^{(t)}}{\partial e^{(t-1)}} = \left( \gamma^{(t-1)} \circ h^{(t-1)} \circ (1 - h^{(t-1)}) \right) \cdot w_e}$$

$$\frac{\partial J^{(t)}}{\partial w_e} \Big|_{t-1} = \left( \delta_2^{(t-1)} \right)^T \cdot \left( e^{(t-1)} \right)^T \quad \begin{matrix} \text{Using results from} \\ (\text{b}) \end{matrix}$$

$$\Rightarrow \boxed{\frac{\partial J^{(t)}}{\partial w_e} \Big|_{t-1} = \left( \gamma^{(t-1)} \circ h^{(t-1)} \circ (1 - h^{(t-1)}) \right)^T \cdot \left( e^{(t-1)} \right)^T}$$

Similarly /

$$\frac{\partial J^{(t)}}{\partial w_h} \Big|_{(t-1)} = \left( \delta_2^{(t-1)} \right)^T \cdot \left( h^{(t-2)} \right)^T \quad \begin{matrix} \text{(Using results} \\ \text{from (b))} \end{matrix}$$

or

$$\boxed{\frac{\partial J^{(t)}}{\partial w_h} \Big|_{(t-1)} = \left( \gamma^{(t-1)} \circ h^{(t-1)} \circ (1 - h^{(t-1)})^T \right)^T \cdot \left( h^{(t-2)} \right)^T}$$

(d)

$$\text{Complexity for } \frac{\partial J^{(t)}}{\partial J} \approx O(|V| \times D_h)$$

$$\text{Complexity for } \frac{\partial J^{(t)}}{\partial e^{(t)}} \approx O(D_h \times d)$$

$$\text{Complexity for } \left. \frac{\partial J^{(t)}}{\partial w_e} \right|_{(t)} \approx O(D_h \times d)$$

$$\text{Complexity for } \left. \frac{\partial J^{(t)}}{\partial w_h} \right|_{(t)} \approx O(D_h \times D_h)$$

Therefore overall complexity

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

(C) If we do backpropagation for  $T$  time steps for one word, the complexity would be:

$$O(VD_h + T(dD_h + D_h^2))$$

And if we hence do backpropagation through time for each of sequence of  $T$  words separately, then it would hence take

$$O(T|V|D_h + T^2(dD_h + D_h^2))$$

Essentially, here we would be passing error signal for each word all the way to  $t=0$  (average  $T$  times). However, we don't need to since all the terms are common except error term. We just need to multiply error term once for each common terms which grows as we move forward. Hence we can reduce the net complexity to

$$O(T(VD_h + dD_h + D_h^2))$$

(f) The term  $O(|V|D_h)$  to compute  $\frac{\partial J^{(t)}}{\partial V}$  is the largest since  $|V| \gg D_h$ . The computational bottleneck is between the hidden and output layers. This term corresponds to the gradient updates of the output layer. The output layer computes the probability distribution of the next words.

### Bonus Question :

The bottleneck is due to the huge vocabulary size. One idea to speed up the computation could be to group the words in a smaller number of classes such that each word belongs to one class only (based on frequency), then we obtain two probability distributions via softmax, one for class prediction and one for word within that class.