

# Midterm Examination

## 1. Tensorflow and Backpropagation

(1)

Using some new notation:

$$z_1 = Wx + b_1$$

$$h = \text{ReLU}(z_1)$$

$$z_2 = Uh + b_2$$

$$\hat{y} = \text{softmax}(z_2)$$

$$J = CE(y, \hat{y})$$

from lectures,

$$\frac{\partial J}{\partial z_2} = \delta_2 = (\hat{y} - y)$$

Backpropagating this error:

$$\boxed{\frac{\partial J}{\partial b_2} = \delta_2}$$

$$\boxed{\frac{\partial J}{\partial U} = \delta_2 h^T}$$

$$\frac{\partial J}{\partial h} = U^T \delta_2$$

$$\frac{\partial J}{\partial z_1} = \delta_1 = \text{ReLU}'(z_1) \odot V^T \delta_2$$

$$\delta_1 = \mathbb{1}\{z_1 > 0\} \odot V^T \delta_2$$

$$\delta_1 = \mathbb{1}\{w_{x_1} + b_1 > 0\} \odot V^T \delta_2$$

$$\boxed{\frac{\partial J}{\partial b_i} = \delta_1}$$

$$\boxed{\frac{\partial J}{\partial w} = \delta_1 x_i^T}$$

(2) They might be required to compute the gradient. for eg. to calculate gradient of sigmoid fun<sup>n</sup>, the o/p of fun<sup>n</sup> is needed.

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

Similarly let  $a = bc$ , then to compute gradient  $\frac{\partial a}{\partial b}$ , the if  $c$  is needed.

$$\frac{\partial a}{\partial b} = c$$

(3) Doing separate softmax will require computing exponents for potentially large numbers, which might result in NaN for the programming language.

On the other hand, CE mode doesn't actually compute exponents of unscaled inputs to compute the loss. Therefore it is more efficient.

## 2. Word2Vec

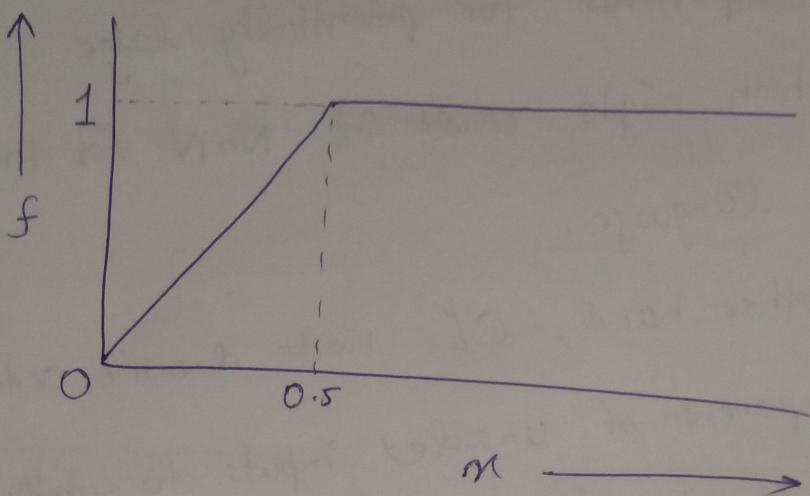
$$\underline{\underline{(1)}} \quad (i) J(\theta) = \frac{1}{2} \sum_{i=1}^w \sum_{j=1}^w f(p_{ij}) (u_i^T v_j - \log p_{ij})^2$$

$$\frac{\partial J(\theta)}{\partial u_i} = \sum_{j=1}^w f(p_{ij}) (u_i^T v_j - \log p_{ij}) v_j$$

$$\frac{\partial J(\theta)}{\partial v_j} = \sum_{i=1}^w f(p_{ij}) (u_i^T v_j - \log p_{ij}) u_i$$

(ii)

$f(n)$  looks like this



(iii)

$f(n)$  should have this form since we want that pairs with high frequency be given higher weightage. We clip it at a maximum value since we don't want most of the weightage given only to those few pairs with very high frequency.

(iv)

Unlike Skipgram/CBOW, Glove provide statistical properties. Also fast to train.

- (2) - Concatenate both vectors of a word  $u_w$  &  $v_w$  to form the representation  $[u_w; v_w]$
- Add the vectors

### 3. Deep NLP in Practice

(1) We will take the distributed representation of each token in the note and concatenate them. Since there are not many tokens, this will not blow up very large.

(2) Adding the word vectors seems a good choice. If a word occurs only few times (e.g. disease name), then averaging would reduce its weights for very large notes, while adding the vectors would preserve such info.

(3) for sigmoid ~~softmax~~, normalization is very important. This is because if a component is too high (or low) then the corresponding linear transformation of NN will result in

(or <sup>very</sup> small)  
a large number. This will result in exploding/vanishing gradient problem for sigmoid & ReLU nonlinearities.

(4) Because RNN learns how the illness of patient has evolved, in a natural way and based on that, it not ~~only~~ learns only the average of past vectors, but where patient vector stands now.

(5) Sigmoid. Since in the softmax case, the presence of one disease would lower the probability of all other diseases. This contradicts our assumption that diseases are not mutually exclusive.

(6) We will now use 300 neuron units instead of 70,000 in the second last layer & ~~one~~  $k$  neurons in a final layer. Here  $k$  is the maximum no. of diseases

under a class (on an average 233).

So, the NN will first predict a class, and then the disease from that class.

$$z_i^{(n-1)} = f_1(w z^{(n-2)} + b_1) \quad \begin{matrix} i \in [800] \\ \text{---} \end{matrix}$$

$$z_i^{(n)} = f_2(v z^{(n-1)} + b_2) \quad \begin{matrix} i \in [K] \\ \text{---} \end{matrix}$$

$$w \in \mathbb{R}^{N \times 300}$$

$$v \in \mathbb{R}^{K \times 300}$$

(7) A bidirectional model would use the future information of patient notes which was not available in current time step to make the prediction at current time step.

(8) Semantic

## 4. LSTMs, GRUs, and Recursive Networks

- (1) (a) Very deep feed-forward NN,  
Recurrent NN & Recursive NNs suffer  
from vanishing gradient problem since all  
of them backpropagate errors to very  
farther in the network. And for networks  
with small errors, it becomes the gradient  
becomes zero as it propagates through the  
network.
- (b) False. Adding more layers will only  
increase the vanishing gradient problem.
- (c) False.  $L_2$  regularizer will push weights of  
neurons towards zero, which will cause  
gradients to become zero.
- (d) True. Clipping gradients will ensure that no  
large gradients are passed to next backprop-  
agated.

(2) (a) False. After multiplication with  $U$  and non-linearities,  $h_t = h_{t-1}$  is generally false.

(b) True. If  $f_t$  is small (or zero), then  $C_{t-1}$  will not have much contribution for calculating  $C_t$  & hence,  $h_t$ .

(c) True. Sigmoid is always b/w 0 & 1.

(d) False. Sigmoid only guarantees that  $o_p$  is b/w 0 & 1. But resulting vector may not sum to 1.

(3)  $h_t \in \mathbb{R}^{d_h \times 1}$  [Assuming real numbers]  
 $x_t \in \mathbb{R}^{d_x \times 1}$   
 $W^{(2)} \in \mathbb{R}^{d_h \times d_h}, W^{(r)} \in \mathbb{R}^{d_h \times d_h},$   
 $U^{(2)} \in \mathbb{R}^{d_h \times d_n}, U^{(s)} \in \mathbb{R}^{d_h \times d_n},$   
 $W \in \mathbb{R}^{d_n \times d_x}, U \in \mathbb{R}^{d_h \times d_n}$

(4)

$$W_2^{(z)} \in \mathbb{R}^{d_n \times d_n}$$

$$U_2^{(z)} \in \mathbb{R}^{d_n \times d_n}$$

$$W_2^{(s)} \in \mathbb{R}^{d_n \times d_n}$$

$$U_2^{(s)} \in \mathbb{R}^{d_n \times d_n}$$

$$W_2 \in \mathbb{R}^{d_n \times d_n}$$

$$U_2 \in \mathbb{R}^{d_n \times d_n}$$

## 5. Hyper-Parameter Tuning

(1)

Bias term acts as the translation & moves the prediction in the direction of original distribution. Adding regularization would prevent this.

(2)

(i) Small weights helps in reducing overfitting. Also, very large weights might cause exploding gradient problem for some deep NNs.

(ii) Initializing weights randomly helps in faster training of the model.

(3) (i) Underfitting

(ii) Make model more complex

(4)

$$\frac{f(x+h) - f(x-h)}{2h} = \sum_{n=0}^{\infty} \frac{f^{(n)}(x+h)(-h)^n}{n!} - \frac{f^{(n)}(x-h)(h)^n}{n!}$$

61  
—  
77