NLP Assignment

Assignment 1
Part 1
Problem 1:
a)
we have $y_w = \begin{cases} 1 \text{ if } w = o \\ 0 \text{ if } w \neq 0 \end{cases} \Rightarrow -\sum\limits_{w \in vocab} y_w \cdot \log(\hat{y}_w) = -y_o \cdot \log(\hat{y}_o)$

$$= -\log(\hat{y}_o)$$

b) $J_{naive-softmax}(v_c, o, U) = -\log p(O = o, C = c)$

$$= -\log \frac{\exp(u_o^T \cdot v_c)}{\sum\limits_{w \in vocab} \exp(u_w^T \cdot v_c)}$$

$$\frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c}\left( -\log \frac{\exp(u_o^T \cdot v_c)}{\sum\limits_{w \in vocab} \exp(u_w^T \cdot v_c)} \right)$$

$$= -\frac{\partial}{\partial v_c} \log\left(\exp(u_o^T \cdot v_c)\right) + \frac{\partial}{\partial v_c} \cdot \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

$$= -\frac{\partial}{\partial v_c} u_o^T \cdot v_c + \frac{\partial}{\partial v_c} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

$$= -u_o + \frac{\partial}{\partial v_c} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

we have:
$$\frac{\partial}{\partial v_c} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right) = \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \frac{\partial}{\partial v_c} \sum\limits_{x \in vocab} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \sum\limits_{x} \frac{\partial}{\partial v_c} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \sum\limits_{x} \exp(u_x^T \cdot v_c) \cdot \frac{\partial}{\partial v_c}(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \sum\limits_{x} \exp(u_x^T \cdot v_c) \cdot u_x$$

$$= \sum\limits_{x} \frac{\exp(u_x^T \cdot v_c)}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot u_x$$

$$= \sum\limits_{x} p(O=x | C=c) \cdot u_x = \sum\limits_{x} \hat{y}_x \cdot u_x$$

So: $\dfrac{\partial J}{\partial v_c} = -u_0 + \sum\limits_{x \in vocab} \hat{y}_x \cdot u_x$

c) $\dfrac{\partial J}{\partial u_w} = \dfrac{\partial}{\partial u_w}\left(-\log \dfrac{\exp(u_0^T \cdot v_c)}{\sum\limits_{w \in vocab} \exp(u_w^T \cdot v_c)}\right)$

$$= -\frac{\partial}{\partial u_w} \log(\exp(u_0^T \cdot v_c)) + \frac{\partial}{\partial u_w} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

$$= -\frac{\partial}{\partial u_w} u_0^T \cdot v_c + \frac{\partial}{\partial u_w} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

In case $w=0$:

$$\frac{\partial J}{\partial u_0} = -\frac{\partial}{\partial u_0} u_0^T \cdot v_c + \frac{\partial}{\partial u_0} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

$$= -v_c + \frac{\partial}{\partial u_0} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right)$$

we have:

$$\frac{\partial}{\partial u_0} \log\left(\sum\limits_{w} \exp(u_w^T \cdot v_c)\right) = \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \frac{\partial}{\partial u_0} \sum\limits_{x \in vocab} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_x \frac{\partial}{\partial u_o} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \sum_x \exp(u_x^T \cdot v_c) \cdot \frac{\partial}{\partial u_o} u_x^T \cdot v_c$$

$$= \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \left( \exp(u_o^T \cdot v_c) \cdot \frac{\partial}{\partial u_o} \cdot u_o^T \cdot v_c + \sum_{x \neq 0} \exp(u_x^T v_c) \cdot \frac{\partial}{\partial u_o} u_x^T \cdot v \right)$$

$$= \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \left( \exp(u_o^T \cdot v_c) \cdot v_c + 0 \right)$$

$$= \frac{\exp(u_o^T \cdot v_c) \cdot v_c}{\sum_w \exp(u_w^T \cdot v_c)} = P(O = o \mid C = c) \cdot v_c = \hat{y}_0 \cdot v_c$$

So: $\dfrac{\partial J}{\partial u_o} = -v_c + \hat{y}_0 \cdot v_c = v_c(1 - \hat{y}_0)$

In case $w \neq 0$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial}{\partial u_w} u_o^T \cdot v_c + \frac{\partial}{\partial u_w} \log\left( \sum_w \exp(u_w^T \cdot v_c) \right)$$

$$= 0 + \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \frac{\partial}{\partial u_w} \sum_{x \in vocab} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \sum_x \frac{\partial}{\partial u_w} \exp(u_x^T \cdot v_c)$$

$$= \frac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \sum_x \exp(u_x^T \cdot v_c) \cdot \frac{\partial}{\partial u_w}(u_x^T \cdot v_c)$$

with each $w \neq 0$: $\dfrac{\partial J}{\partial u_w} = \dfrac{1}{\sum_w \exp(u_w^T \cdot v_c)} \cdot \left( \dfrac{\partial}{\partial u_w}(u_x^T) \right)$

$$= \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \left( \exp(u_w^T \cdot v_c) \cdot \frac{\partial}{\partial u_w}(u_w^T \cdot v_c) \right.$$

$$\left. + \sum\limits_{x \neq w} \exp(u_x^T \cdot v_c) \cdot \frac{\partial}{\partial u_w}(u_x^T \cdot v_c) \right)$$

$$= \frac{1}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot \left( \exp(u_w^T \cdot v_c) \cdot v_c + 0 \right)$$

$$= \frac{\exp(u_w^T \cdot v_c)}{\sum\limits_{w} \exp(u_w^T \cdot v_c)} \cdot v_c = P(O=w \mid C=c) \cdot v_c$$

$$= \hat{y}_w \cdot v_c$$

with $w \neq 0$

d)

$$\frac{d}{dx}\sigma(x) = \frac{d}{dx} \frac{1}{1+e^{-x}} = \frac{d}{dx}(1+e^{-x})^{-1}$$

$$= -(1+e^{-x})^{-2} \cdot \frac{d}{dx}(1+e^{-x})$$

$$= -(1+e^{-x})^{-2} \cdot \left\{ \frac{d}{dx} e^{-x} \right. = (1+e^{-x})^{-2} \cdot e^{-x} \cdot \frac{d}{dx}(-x)$$

$$= (1+e^{-x})^{-2} \cdot e^{-x}$$

$$= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}}$$

$$= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}+1-1}{1+e^{-x}}$$

$$= \frac{1}{1+e^{-x}} \cdot \left( 1 - \frac{1}{1+e^{-x}} \right)$$

$$= \sigma(x)(1-\sigma(x))$$

e) $J_{neg-sample}(v_c, o, U) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1}^{k} \log(\sigma(-u_k^T \cdot v_c))$

$$\frac{\partial J}{\partial v_c} = \underbrace{-\frac{\partial}{\partial v_c} \log(\sigma(u_o^T \cdot v_c))}_{(1)} - \underbrace{\frac{\partial}{\partial v_c} \sum_k \log(\sigma(-u_k^T \cdot v_c))}_{(2)}$$

$(1) = -\frac{1}{\sigma(u_o^T \cdot v_c)} \cdot \frac{\partial}{\partial v_c} \sigma(u_o^T \cdot v_c)$

$= -\frac{1}{\sigma(u_o^T \cdot v_c)} \cdot \sigma(u_o^T \cdot v_c) \cdot (1 - \sigma(u_o^T \cdot v_c)) \cdot \frac{\partial}{\partial v_c} u_o^T \cdot v_c$

$= (\sigma(u_o^T \cdot v_c) - 1) \cdot u_o$

$(2) = \sum_k \frac{\partial}{\partial v_c} \log(\sigma(-u_k^T \cdot v_c))$

$= \sum_k \frac{1}{\sigma(-u_k^T \cdot v_c)} \cdot \frac{\partial}{\partial v_c} \sigma(-u_k^T \cdot v_c)$

$= \sum_k \frac{1}{\sigma(-u_k^T \cdot v_c)} \cdot \sigma(-u_k^T \cdot v_c)(1 - \sigma(-u_k^T \cdot v_c)) \frac{\partial}{\partial v_c}(-u_k^T \cdot v_c)$

$= \sum_k (\sigma(-u_k^T \cdot v_c) - 1) \cdot u_k$

So: $\frac{\partial J}{\partial v_c} = (\sigma(u_o^T \cdot v_c) - 1) \cdot u_o - \sum_{k=1}^{k} (\sigma(-u_k^T \cdot v_c) - 1) \cdot u_k$

$$\frac{\partial J}{\partial u_w} = \frac{-\partial}{\partial u_w} \log(\sigma(u_o^T \cdot v_c)) - \frac{\partial}{\partial u_w} \sum_k \log(\sigma(-u_k^T \cdot v_c))$$

$$= -(1-\sigma(u_o^T \cdot v_c)) \frac{\partial}{\partial u_w} u_o^T v_c - \sum_k (1-\sigma(-u_k^T \cdot v_c)) \frac{\partial \sigma(-u_k^T \cdot v_c)}{\partial u_w}$$

In case $w=0$, we have:

$$\frac{\partial J}{\partial u_o} = (\sigma(u_o^T \cdot v_c) - 1) \cdot v_c - 0 \quad // \text{ this is because } 0 \notin \{u_1 ... u_k\}$$
$$\text{so the derivative } = 0$$

$$= (\sigma(u_o^T \cdot v_c) - 1) \cdot v_c$$

In case $w \neq 0$, we have:

$$\frac{\partial J}{\partial u_w} = 0 - \sum_k (1-\sigma(-u_k^T \cdot v_c)) \cdot \frac{\partial}{\partial u_w} (-u_k^T \cdot v_c)$$

$$= (1-\sigma(-u_k^T \cdot v_c)) \cdot v_c$$

// this because there is only one $k=w$ so the
derivative of other $= 0$

Conclusion:
- This loss function is much better to compute more
than the naive - softmax loss because it doesn't go
through all the word in the vocabulary therefore,
the computation is less expensive.

**Problem 2:**

a)

— Neural window-based models can be parallelized, but RNN models cannot

b)

— Use a network with fewer layers

— Increase L2 regularization weight

c) — $\log (0, 2)$

d)

— True

e)

— True

**Problem 4:**

i) 1. Number of outputs: n

// because each output will correspond to an input at each time step.

2. $\hat{y}^{(t)}$ is the probability distribution over 4 categories including: person, organization, location, none

3. Each input will be a word in the sentence and will produce an output ~~at each~~ correspond to predicted category at each time step

ii) 1. Number of outputs: arbitrary

// because we don't know how many words the model will generate

2. $\hat{y}^{(t)}$ is the probability distribution over all words in vocabulary

3. Each input will be ~~the output of~~ the previous output and will produce a new output which is the next predicted word for the sentence at each time step

## Problem 3:

a) Yes, graph-based dependency can produce non-projective ~~parsing these~~ dependency trees while transition-based can't

b) For example, in case $i = 5$, the score will be
$$s(i,i) = h_i^T \cdot h_i$$
The result of this product will always be positive and higher than the scores of other edges

⇒ This can cause ~~a~~ a problem that words can get linked to themselves

⇒ Incorrectly predict dependencies that involve a word and itself.

c) The drawback is that graph-based is slower and require more computational resources ~~so~~ since it has to compute all the scores for all possible dependencies between all pairs of words in the sentence while transition-based ~~only~~ just simply build a tree from left to right or right to left incrementally