

# CS 224N Assignment #2

## 1. Written: Understanding word2vec

(a) Note  $\vec{y}$  is a one-hot vector with one at index  $w=0$

$$\begin{aligned}
 - \sum_{w \in V} y_w \log \hat{y}_w &= - (y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + \dots + y_{|V|} \log \hat{y}_{|V|}) \\
 &= - \log \hat{y}_0
 \end{aligned}$$

(b)  $J_{\text{native-softmax}}(V_c, o, U) = -\log P(D=o | C=c)$

$$= -\log \frac{e^{u_o^T V_c}}{\sum_{w \in \text{Vocab}} e^{u_w^T V_c}}$$

$$= -u_o^T V_c + \log \sum_{w \in \text{Vocab}} e^{u_w^T V_c}$$

$$\frac{\partial J_{\text{native}}}{\partial V_c} = -\frac{\partial}{\partial V_c} u_o^T V_c + \frac{\partial}{\partial V_c} \log \sum_{w \in \text{Vocab}} e^{u_w^T V_c}$$

$$= -u_o + \frac{1}{\sum_{w \in \text{Vocab}} e^{u_w^T V_c}} \cdot \sum_{w \in \text{Vocab}} u_w \cdot e^{u_w^T V_c}$$

$$= -u_o + \sum_{w \in \text{Vocab}} \frac{e^{u_w^T V_c}}{\sum_{w \in \text{Vocab}} e^{u_w^T V_c}} \cdot u_w$$

$$= -u_o + \sum_{w \in \text{Vocab}} P(u_w | V_c) \cdot u_w$$

$$= -u_o + \sum_{w \in \text{Vocab}} \hat{y}_w \cdot u_w$$

(c) If  $w = 0$

$$\begin{aligned}
 \frac{\partial J_{\text{native}}}{\partial u_0} &= -\frac{\partial}{\partial u_0} u_0^T v_c + \frac{\partial}{\partial u_0} \log \sum_{w \in \text{Vocab}} e^{u_w^T v_c} \\
 &= -v_c + \frac{1}{\sum_{w \in \text{Vocab}} e^{u_w^T v_c}} \cdot \frac{\partial}{\partial u_0} \sum_{w \in \text{Vocab}} e^{u_w^T v_c} \\
 &= -v_c + \frac{1}{\sum_{w \in \text{Vocab}} e^{u_w^T v_c}} \cdot e^{u_0^T v_c} \cdot v_c \\
 &= v_c (\hat{y}_0 - 1)
 \end{aligned}$$

else (i.e.,  $w \neq 0$ )

$$\begin{aligned}
 \frac{\partial J_{\text{native}}}{\partial u_w} &= -\frac{\partial}{\partial u_w} u_0^T v_c + \frac{\partial}{\partial u_w} \log \sum_{x \in \text{Vocab}} e^{u_x^T v_c} \\
 &= 0 + \frac{1}{\sum_{x \in \text{Vocab}} e^{u_x^T v_c}} \frac{\partial}{\partial u_w} \sum_{x \in \text{Vocab}} e^{u_x^T v_c} \\
 &= \frac{1}{\sum_{x \in \text{Vocab}} e^{u_x^T v_c}} e^{u_w^T v_c} \cdot v_c \\
 &= v_c \hat{y}_w
 \end{aligned}$$

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

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

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

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

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

$$(e) \quad J_{\text{neg}} = -\log(\sigma(u_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T V_c))$$

$$= -\log \frac{1}{1 + e^{-u_0^T V_c}} - \sum_{k=1}^K \log \frac{1}{1 + e^{u_k^T V_c}}$$

compute partial derivative for  $V_c$

$$\frac{\partial J_{\text{neg}}}{\partial V_c} = - \frac{\partial}{\partial V_c} \log \frac{1}{1 + e^{-u_0^T V_c}} - \sum_{k=1}^K \frac{\partial}{\partial V_c} \log \frac{1}{1 + e^{u_k^T V_c}}$$

$$= - \frac{1}{\sigma(u_0^T V_c)} \cdot \sigma(u_0^T V_c) \cdot [1 - \sigma(u_0^T V_c)] \cdot u_0$$

$$- \sum_{k=1}^K \frac{1}{\sigma(-u_k^T V_c)} \cdot \sigma(-u_k^T V_c) \cdot [1 - \sigma(-u_k^T V_c)] \cdot$$

$$(-u_k)$$

$$= -u_0 [1 - \sigma(u_0^T V_c)] + \sum_{k=1}^K u_k [1 - \sigma(-u_k^T V_c)]$$

Compute partial derivative for  $u_0$

$$\begin{aligned}\frac{\partial J_{\text{neg}}}{\partial u_0} &= - \frac{1}{\sigma(u_0^T v_c)} \cdot \sigma(u_0^T v_c) \cdot [1 - \sigma(u_0^T v_c)] \cdot v_c \\ &\quad - \sum_{k=1}^k \frac{\partial}{\partial u_0} \log(\sigma(-u_k^T v_c)) \\ &= -v_c [1 - \sigma(u_0^T v_c)]\end{aligned}$$

Compute partial derivative for  $u_k$

$$\begin{aligned}\frac{\partial J_{\text{neg}}}{\partial u_k} &= - \frac{\partial}{\partial u_k} \log(\sigma(u_0^T v_c)) - \sum_{x=1}^k \frac{\partial}{\partial u_k} \log(\sigma(-u_x^T v_c)) \\ &= - \frac{\partial}{\partial u_k} \log(\sigma(-u_k^T v_c)) \\ &= - \frac{1}{\sigma(-u_k^T v_c)} \cdot \sigma(-u_k^T v_c) \cdot [1 - \sigma(-u_k^T v_c)] \cdot (-v_c) \\ &= v_c [1 - \sigma(-u_k^T v_c)]\end{aligned}$$

For the naive softmax loss function, we will have to use all the word vectors in the vocabulary to normalize the probabilities. In contrast, with negative sampling, we just need to take  $k$  word vectors into account.

$$(f) \quad J_{\text{skip-gram}}(V_c, W_{t-m}, \dots, W_{t+m}, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(V_c, W_{t+j}, U)$$

$$(i) \quad \frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(V_c, W_{t+j}, U)}{\partial U}$$

$$(ii) \quad \frac{\partial J_{\text{skip-gram}}}{\partial V_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(V_c, W_{t+j}, U)}{\partial V_c}$$

$$(iii) \quad \frac{\partial J_{\text{skip-gram}}}{\partial V_w} = 0$$