## 1. Written : Understanding Word2Vec

**(a)**

$$y_w = \begin{cases} 1 & (w=o) \\ 0 & (w \neq o) \end{cases}$$

$$- \sum_{w \in Vocab} y_w \log \hat{y}_w = - y_o \log \hat{y}_o - \sum_{w \neq o} y_w \log \hat{y}_w = - \log \hat{y}_o$$

$$J_{naive-softmax} = - \log \frac{\exp(u_o^T v_c)}{\sum_{w \in vocab} \exp(u_w^T v_c)} = - \log \hat{y}_o$$

**(b)**

$$\frac{\partial J_{naive-softmax}}{\partial v_c} = \frac{\partial}{\partial v_c} \left( - \log \frac{\exp(u_o^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)} \right) = \frac{\partial}{\partial v_c} \left( - \log \exp(u_o^T v_c) + \log \sum_{w \in Vocab} \exp(u_w^T v_c) \right)$$

$$= - \frac{1}{\exp(u_o^T v_c)} \cdot \exp(u_o^T v_c) \cdot u_o + \frac{1}{\sum_{w \in Vocab} \exp(u_w^T v_c)} \cdot \left( \sum_{w \in Vocab} \exp(u_w^T v_c) \cdot u_w \right)$$

$$= - u_o + \sum_{w \in Vocab} \left( \frac{\exp(u_w^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)} \cdot u_w \right) = - u_o + \sum_{w \in Vocab} \hat{y}_w \cdot u_w$$

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

$$= \sum_{w \in Vocab} \underbrace{(\hat{y}_w - y_w)}_{scalar} \cdot \underbrace{u_w}_{vector (d \times 1)} \in \mathbb{R}^{d \times 1}$$

$$\left( U = [u_1, u_2, \cdots, u_o, \cdots u_{|Vocab|}] \in \mathbb{R}^{d \times |Vocab|} \right.$$
$$\left. y \in \mathbb{R}^{|Vocab| \times 1}, \ \hat{y} \in \mathbb{R}^{|Vocab| \times 1} \right.$$

$$\therefore \frac{\partial J_{naive-softmax}}{\partial v_c} = U \cdot (\hat{y} - y)$$

**(c)**

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

i) $w = o$

$$\frac{\partial}{\partial u_w} \left( - u_o^T v_c + \log \sum_w \exp(u_w^T v_c) \right) = - v_c + \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \cdot v_c = - v_c + \hat{y}_w \cdot v_c = (\hat{y}_w - 1) \cdot v_c$$

ii) $w \neq o$

$$\frac{\partial}{\partial u_w} \left( - u_o^T v_c + \log \sum_w \exp(u_w^T v_c) \right) = 0 + \frac{\exp(u_w^T v_c)}{\sum \exp(u_w^T v_c)} \cdot v_c = \hat{y}_w \cdot v_c$$

$$(\because) \ \frac{\partial J_{naive-softmax}}{\partial u_w} = (\hat{y}_w - y_w) v_c \begin{cases} (\hat{y}_w - 1) \cdot v_c & (w = o) \\ \hat{y}_w \cdot v_c & (w \neq o) \end{cases}$$

**(d)**

$$\frac{\partial J_{naive-softmax}}{\partial U} = \left[ \frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \cdots, \frac{\partial J}{\partial u_{|Vocab|}} \right] = \left[ (\hat{y}_1 - y_1) \cdot v_c, (\hat{y}_2 - y_2) \cdot v_c, \cdots, (\hat{y}_{|Vocab|} - y_{|Vocab|}) \cdot v_c \right]$$

$$= v_c \cdot (\hat{y} - y)^T$$

**(e)**

$$J_{neg}(v_c, o, U) = - \log \sigma(u_o^T v_c) - \sum_{s=1}^{K} \log \sigma(- u_{w_s}^T v_c)$$

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

$$\frac{\partial J}{\partial v_c} = - \frac{1}{\sigma(u_o^T v_c)} \cdot \sigma(u_o^T v_c) \cdot (1 - \sigma(u_o^T v_c)) \cdot u_o - \sum_{s=1}^{K} \frac{1}{\sigma(-u_{w_s}^T v_c)} \cdot \sigma(-u_{w_s}^T v_c)(1 - \sigma(-u_{w_s}^T v_c)) \cdot (-u_{w_s})$$

$$= (\sigma(u_o^T v_c) - 1) \cdot u_o + \sum_{s=1}^{K} (1 - \sigma(-u_{w_s}^T v_c)) \cdot u_{w_s} = (\sigma(u_o^T v_c) - 1) \cdot u_o + \sum_{s=1}^{K} \sigma(u_{w_s}^T v_c) \cdot u_{w_s}$$

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

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

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

$$(\because) \begin{cases} \frac{\partial J}{\partial v_c} = (\sigma(u_o^T v_c) - 1) \cdot u_o + \sum_{s=1}^{K} \sigma(u_{w_s}^T v_c) \cdot u_{w_s} \\ \frac{\partial J}{\partial u_o} = (\sigma(u_o^T v_c) - 1) \cdot v_c \\ \frac{\partial J}{\partial u_{w_s}} = \sigma(u_{w_s}^T v_c) \cdot u_{w_s} \end{cases}$$

(ii) Negative sampling only computes $K$ terms,
which is much smaller than total vocabulary.

**(f)**

(i) $\frac{\partial J_{skip}}{\partial U} = - \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U}$

(ii) $\frac{\partial J_{skip}}{\partial v_c} = - \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}$

(iii) $\frac{\partial J_{skip}}{\partial v_w} (w \neq c) = 0$ (No term of $v_w$ → Derivatives will be zero.)