

$$a) - \sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -y_{w'} \log(\hat{y}_{w'}) = -\log(\hat{y}_{w'})$$

\* We know that all of  $y_w | w \in \text{vocab}$  are zeros except one of them. Thus, we can omit all words with  $y_w = 0$  and keep  $w'$  which  $y_{w'} = 1$

$$-\log(\hat{y}_{w'}) = -\log(\hat{y}_0) = -\log(P(O=0|C=c)) = J_{\text{naive softmax}}$$

\* We know that  $\hat{y}_0$  is scalar. Thus it can be written instead of  $\hat{y}_{w'}$ .

$$\begin{aligned} b)(i) \quad \frac{\partial J_{\text{naive-softmax}}}{\partial v_c} &= - \frac{\partial \log(P(O=0|C=c))}{\partial v_c} \\ &= - \frac{\partial \log\left(\frac{\exp(u_0^T \cdot v_c)}{\sum \exp(u_w^T \cdot v_c)}\right)}{\partial v_c} = - \frac{\partial \log(\exp(u_0^T \cdot v_c))}{\partial v_c} \\ &\quad + \frac{\partial \log(\sum \exp(u_w^T \cdot v_c))}{\partial v_c} = - \frac{1}{\exp(u_0^T \cdot v_c)} \cdot \frac{\partial \exp(u_0^T \cdot v_c)}{\partial v_c} \\ &\quad + \frac{1}{\sum \exp(u_w^T \cdot v_c)} \cdot \frac{\partial \sum \exp(u_w^T \cdot v_c)}{\partial v_c} = - \frac{\exp(u_0^T \cdot v_c)}{\exp(u_0^T \cdot v_c)} \cdot u_0 \\ &\quad + \frac{\sum u_w \exp(u_w^T \cdot v_c)}{\sum \exp(u_w^T \cdot v_c)} = -u_0 + \sum_{w \in \text{vocab}} u_w \left( \frac{\exp(u_w^T \cdot v_c)}{\sum_{w' \in \text{vocab}} \exp(u_{w'}^T \cdot v_c)} \right) \\ &= -u_0 + \sum_{w \in \text{vocab}} u_w P(O=w|C=c) = -u_0 + U \cdot \hat{y} \end{aligned}$$

$u_0 \in \mathbb{R}^{d \times 1}, \quad U \in \mathbb{R}^{d \times \text{vocab-size}}, \quad \hat{y} \in \mathbb{R}^{\text{vocab-size} \times 1}$

$$(ii) \quad -u_o + U \cdot \hat{y} = 0 \Rightarrow u_o = U \cdot \hat{y}$$

$\downarrow$  outside context word       $\rightarrow$  predicted context word

\* It means that, gradient will be zero, if the predicted word is equal to outside context word.

(iii) \* The gradient that we found is the error of our model in predicting the correct outside context word. Thus, by subtracting this error from  $v_c$  we are moving the vector which is representing  $v_c$  to a new position in  $d$ -dimensional space in order to force our model to predict words with minimum difference with the actual outside context word.

(iv)

$L_2$  normalization is forcing vectors to have equal length.

This normalization takes away useful information when some features of a specific word is dividable by the same feature of another word.

For example:

$$\text{royalty}_{\text{king}} = \alpha \cdot \text{royalty}_{\text{prince}}$$

$L_2$  normalization makes these two words look equal in the term of royalty.



c)  $w=0$ :

$$\begin{aligned}
 \frac{\partial \mathcal{J}_{\text{naive-softmax}}}{\partial u_w} &= \frac{-\partial \log\left(\frac{\exp(u_o^T \cdot v_c)}{\sum \exp(u_x^T \cdot v_c)}\right)}{\partial u_w} \\
 &= -\frac{\partial \log(\exp(u_o^T \cdot v_c))}{\partial u_w} + \frac{\partial \log(\sum \exp(u_x^T \cdot v_c))}{\partial u_w} \\
 &= -\frac{1}{\exp(u_o^T \cdot v_c)} \cdot \frac{\exp(u_o^T \cdot v_c)}{\partial u_w} + \frac{1}{\sum_{x \in \text{vocab}} \exp(u_x^T \cdot v_c)} \cdot \frac{\sum \exp(u_x^T \cdot v_c)}{\partial u_w} \\
 &= -v_c + \frac{\exp(u_w^T \cdot v_c)}{\sum_{x \in \text{vocab}} \exp(u_x^T \cdot v_c)} \cdot v_c^T = -v_c + \hat{y}_w \cdot v_c^T
 \end{aligned}$$

$$v_c \in \mathbb{R}^{d \times 1} \quad \hat{y}_w = \mathbb{R}^{1 \times 1}$$

$w \neq 0$ :

$$\begin{aligned}
 \frac{\partial \mathcal{J}_{\text{naive-softmax}}}{\partial u_w} &= \frac{-\partial \log\left(\frac{\exp(u_o^T \cdot v_c)}{\sum \exp(u_x^T \cdot v_c)}\right)}{\partial u_w} \\
 &= -\frac{\partial \log(\exp(u_o^T \cdot v_c))}{\partial u_w} + \frac{\partial \log(\sum_{x \in \text{vocab}} \exp(u_x^T \cdot v_c))}{\partial u_w} \\
 &= 0 + \frac{1}{\sum_{x \in \text{vocab}} \exp(u_x^T \cdot v_c)} \cdot \frac{\partial \sum \exp(u_x^T \cdot v_c)}{\partial u_w} = \frac{\exp(u_w^T \cdot v_c)}{\sum_{x \in \text{vocab}} \exp(u_x^T \cdot v_c)} \cdot v_c^T \\
 &= \hat{y}_w \cdot v_c^T
 \end{aligned}$$

$$v_c \in \mathbb{R}^{d \times 1} \quad \hat{y}_w \in \mathbb{R}^{1 \times 1}$$

$$d) \quad \frac{\partial J_{\text{naive-softmax}}}{\partial u_w} \in \mathbb{R}^{d \times 1}$$

$$\rightarrow \frac{\partial J_{\text{naive-softmax}}}{\partial u} = \left[ \frac{\partial J_{\text{naive-softmax}}}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_{\text{vocab-len}}} \right]$$

$$e) \quad f(x) = \begin{cases} x & x \geq 0 \\ \alpha x & x < 0 \end{cases} \Rightarrow f'(x) = \begin{cases} 1 & x > 0 \\ \alpha & x < 0 \end{cases}$$

\* This function is not differentiable at  $x=0$ .

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

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

g) (i.)

$$\begin{aligned}
 \frac{\partial J_{\text{neg-sample}}}{\partial v_c} &= \frac{\partial \log(\sigma(u_o^T v_c))}{\partial v_c} - \frac{\partial \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c))}{\partial v_c} \\
 &= -\frac{1}{\sigma(u_o^T v_c)} \cdot \frac{\partial \sigma(u_o^T v_c)}{\partial v_c} - \sum_{s=1}^K \frac{1}{\sigma(-u_{w_s}^T v_c)} \cdot \frac{\partial \sigma(-u_{w_s}^T v_c)}{\partial v_c} \\
 &= -\frac{\sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c))}{\sigma(u_o^T v_c)} \cdot \frac{\partial u_o^T v_c}{\partial v_c} \\
 &\quad - \sum_{s=1}^K \frac{\sigma(-u_{w_s}^T v_c)(1 - \sigma(-u_{w_s}^T v_c))}{\sigma(-u_{w_s}^T v_c)} \cdot \frac{\partial (-u_{w_s}^T v_c)}{\partial v_c} \\
 &= -u_o(1 - \sigma(u_o^T v_c)) + \sum_{s=1}^K u_{w_s}(1 - \sigma(-u_{w_s}^T v_c))
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial J_{\text{neg-sample}}}{\partial u_o} &= -\frac{\partial \log(\sigma(u_o^T v_c))}{\partial u_o} - \frac{\partial \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c))}{\partial u_o} \\
 &= -\frac{\sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c))}{\sigma(u_o^T v_c)} \cdot \frac{\partial (u_o^T v_c)}{\partial u_o} - 0
 \end{aligned}$$

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

$$\begin{aligned}
 \frac{\partial J_{\text{neg-sample}}}{\partial u_{w_x}} &= -\frac{\partial \log(\sigma(u_o^T v_c))}{\partial u_{w_x}} - \frac{\partial \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c))}{\partial u_{w_x}} \\
 &= 0 - \frac{\sigma(-u_{w_x}^T v_c)(1 - \sigma(-u_{w_x}^T v_c))}{\sigma(-u_{w_x}^T v_c)} \cdot \frac{\partial (-u_{w_x}^T v_c)}{\partial u_{w_x}}
 \end{aligned}$$

$$= (1 - \sigma(-u_{w_x}^T v_c)) v_c^T$$



$$(ii) \quad U_{0, \{w_1, \dots, w_k\}}^T \cdot v_c = \begin{bmatrix} u_0^T \cdot v_c \\ -u_{w_1}^T \cdot v_c \\ \vdots \\ -u_{w_k}^T \cdot v_c \end{bmatrix} \Rightarrow 1 - \sigma(U_{0, \{w_1, \dots, w_k\}}^T \cdot v_c) = \begin{bmatrix} 1 - \sigma(u_0^T \cdot v_c) \\ 1 - \sigma(u_{w_1}^T \cdot v_c) \\ \vdots \\ 1 - \sigma(u_{w_k}^T \cdot v_c) \end{bmatrix}$$

(iii) We don't need to iterate over the vocabulary.

h) All the terms in  $\frac{\partial \tilde{J}_{\text{neg-sample}}}{\partial u_{w_x}}$  will be zero, except

the ones with  $u_{w_x} = u_{w_s}$

$$\frac{\partial \tilde{J}_{\text{neg-sample}}}{\partial u_{w_x}} = - \frac{\sum_{\substack{1 \leq s \leq k \\ w_s = w_x}} \log(\sigma(-u_{w_s}^T \cdot v_c))}{\partial u_{w_x}}$$

$$= - \sum_{\substack{1 \leq s \leq k \\ w_s = w_x}} \frac{\sigma(-u_{w_s}^T \cdot v_c)(1 - \sigma(-u_{w_s}^T \cdot v_c))}{\partial u_{w_x}} \cdot \frac{\partial (-u_{w_s}^T \cdot v_c)}{\partial u_{w_x}}$$

$$= \sum_{\substack{1 \leq s \leq k \\ w_s = w_x}} (1 - \sigma(u_{w_s}^T \cdot v_c)) v_c^T$$

Subject: \_\_\_\_\_

Date \_\_\_\_\_

i)

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

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

$$(iii) \quad \frac{\partial J_{\text{skip-gram}}}{\partial u_w} \text{ when } w \neq c = 0$$